

# **Data Analysis III**

Beate Heinemann

*UC Berkeley and Lawrence Berkeley National Laboratory*

Hadron Collider Physics Summer School, Fermilab, August 2008

# Outline

- Lecture I:
  - Measuring a cross section
    - focus on acceptance
- Lecture II:
  - Measuring a property of a known particle
- Lecture III:
  - Searching for a new particle
    - focus on backgrounds

# Search for New Particles: Experimentally

Number of observed  
events: counted

Background:  
Measured from data /  
calculated from theory

$$\sigma_{\text{NP}} = \frac{N_{\text{obs}} - N_{\text{BG}}}{\int L dt \cdot \epsilon}$$

Cross section  $\sigma$

Luminosity:  
Determined by accelerator,  
trigger prescale, ...

Efficiency:  
optimized by  
experimentalist

- Exactly like with measuring the cross section...

# But we need to observe first!

- When we don't know if a particle exists our first question is: "Does it exist?"
- => significance of signal
  - I.e. how consistent is the number of observed events with the number of background events?
  - Background expectation:  $N_{BG}$ 
    - Expect it to fluctuate statistically by  $\delta N_{BG} \sim \sqrt{N_{BG}}$
  - Signal expectation:  $N_{Signal}$
  - Statistical Significance:  $N_{Signal}/\delta N_{BG} \sim N_{Signal} / \sqrt{N_{BG}}$ 
    - Often called  $S/\sqrt{B}$

in Gaussian limit

	evidence	observation
significance	$3\sigma$	$5\sigma$
Probability of stat. fluctuation	0.3%	$5.7 \times 10^{-8}$

In real life: systematic uncertainties also contribute to  $\delta N_{BG}$



# Search analyses

- Primary focus is background estimate
  - Determines whether or not an observation can be made
  - Cuts for background reduction studied often using benchmark New Physics scenario
    - Also model-independent analyses..
- Secondary focus is acceptance/efficiency determination: required only
  - when putting an upper limit on a cross section
  - when measuring the cross section of the observed new particle
    - Need to know what it is though
    - Or quote cross section for some effective cuts

# Example Analyses

- SUSY:
  - Squarks/gluinos  $\rightarrow$  jets +  $E_T$  (+leptons)
- Higgs:
  - Higgs  $\rightarrow$  WW
- High Mass Resonance
  - $Z'$  (if I have time)

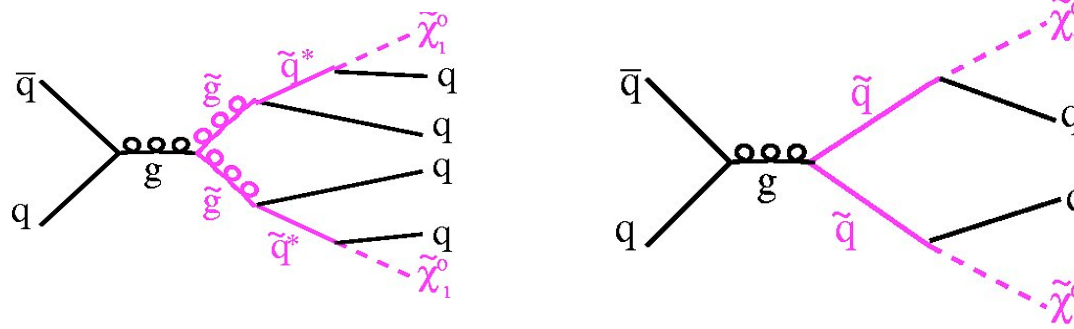
# Backgrounds

- Ideally you get the backgrounds to be small
  - The smaller they are the less well you need to know them
- Estimates based on
  - Data only
    - E.g. lepton fake rates
  - Monte Carlo only
    - For well known electroweak processes
  - Monte Carlo / Data hybrid
    - For e.g. W/Z+jets or W/Z+b-jets

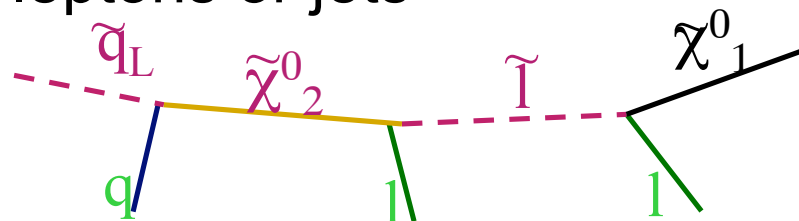
**Squarks/Gluinos →  
Jets + MEt (+ leptons)**

# SUSY at the LHC

- Cross section **much** higher than at Tevatron, e.g.
  - for  $m(\tilde{g})=400$  GeV:  $\sigma_{\text{LHC}}(\tilde{g}\tilde{g})/\sigma_{\text{Tevatron}}(\tilde{g}\tilde{g})\approx 20,000$
  - for  $m(\tilde{q})=400$  GeV:  $\sigma_{\text{LHC}}(\tilde{g}\tilde{g})/\sigma_{\text{Tevatron}}(\tilde{g}\tilde{g})\approx 1,000$ 
    - Since there are a lot more gluons at the LHC (lower x)

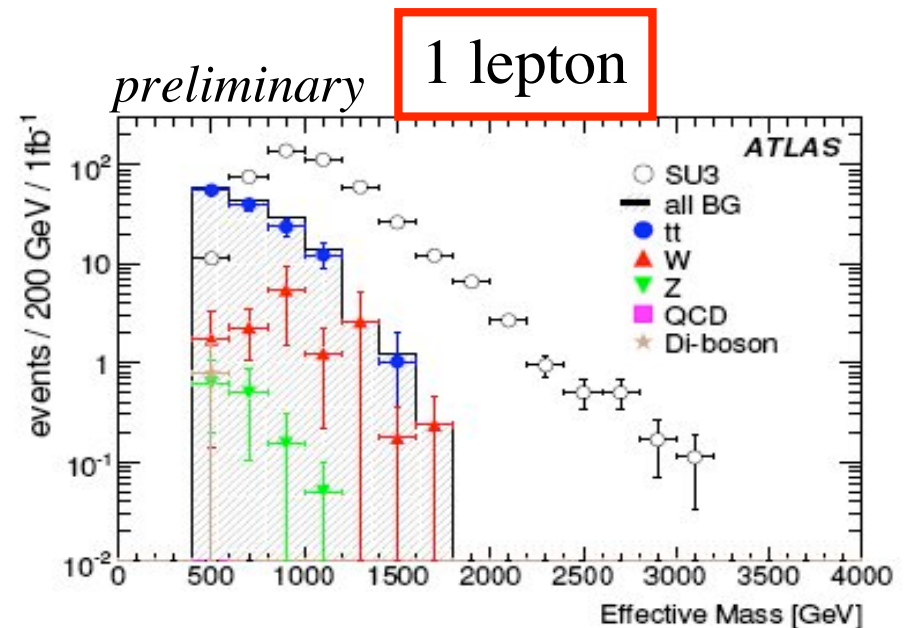
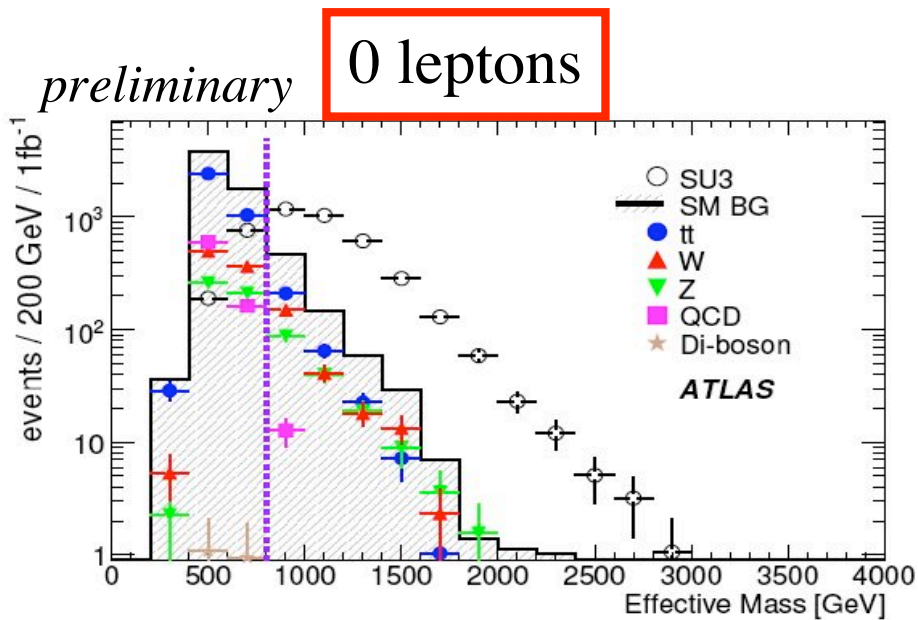
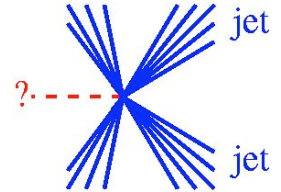


- At higher masses more phase space to decay in cascades
  - Results in additional leptons or jets



# SUSY at the LHC

- Example:  $m(\tilde{q}) \sim 600$  GeV,  $m(\tilde{g}) \sim 700$  GeV
- Require 4 jets, large missing  $E_T$  and 0 or 1 lepton

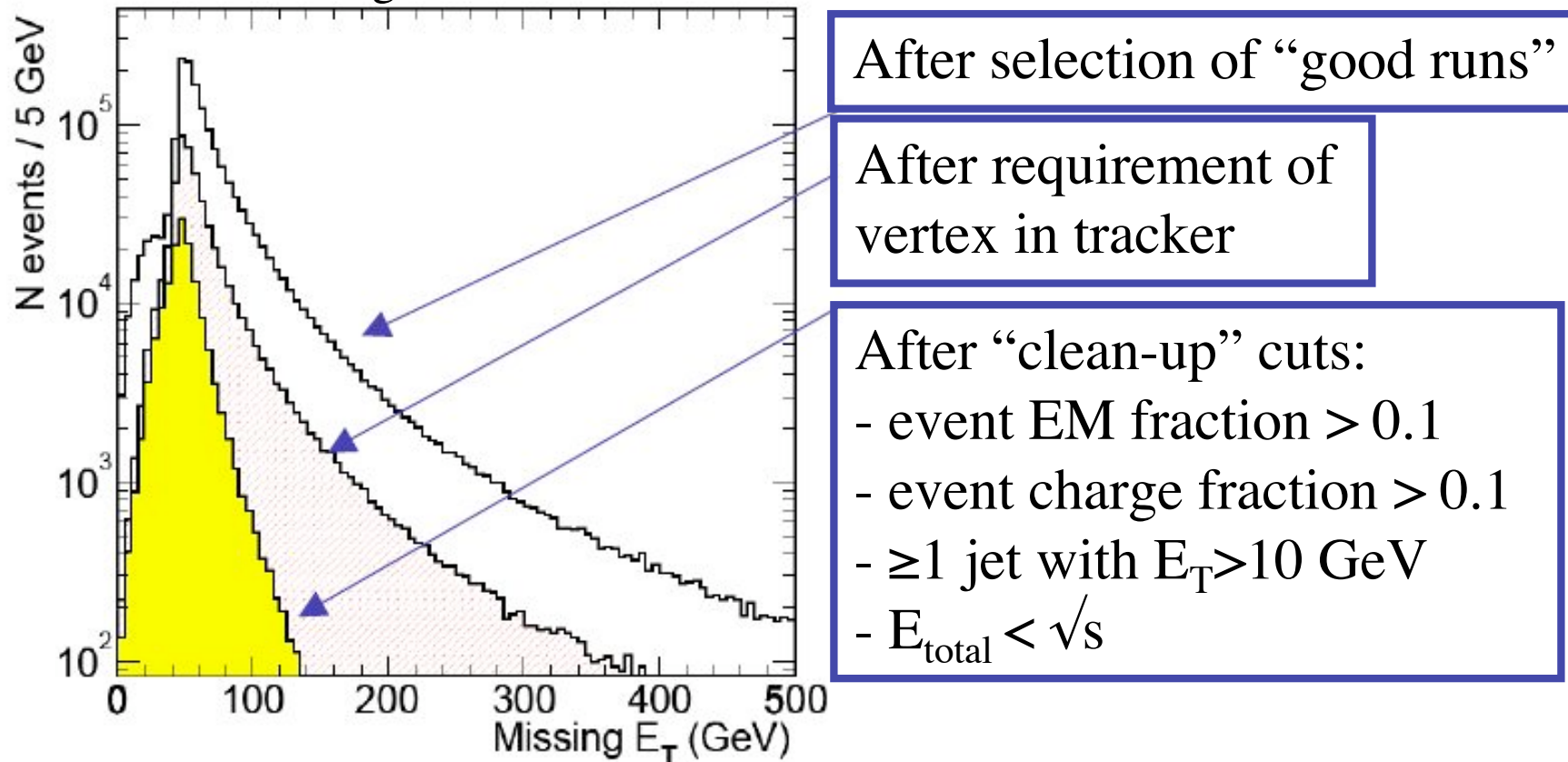


- Similar and great (!) sensitivity in both modes
- Main backgrounds: top, W/Z+jets, QCD multi-jet

But how do we know the backgrounds!?!

# Instrumental Backgrounds

from Avi Yagil



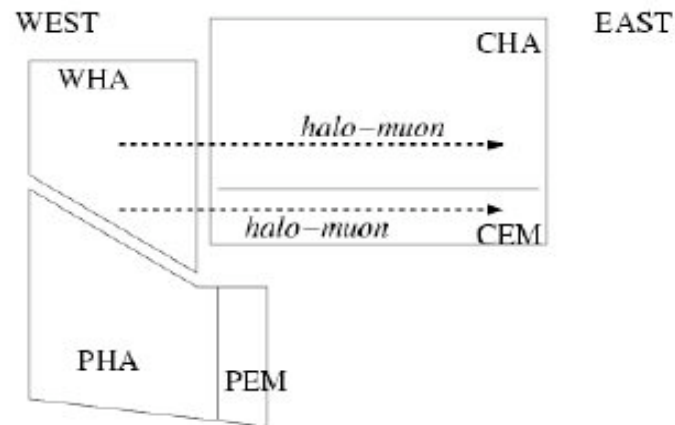
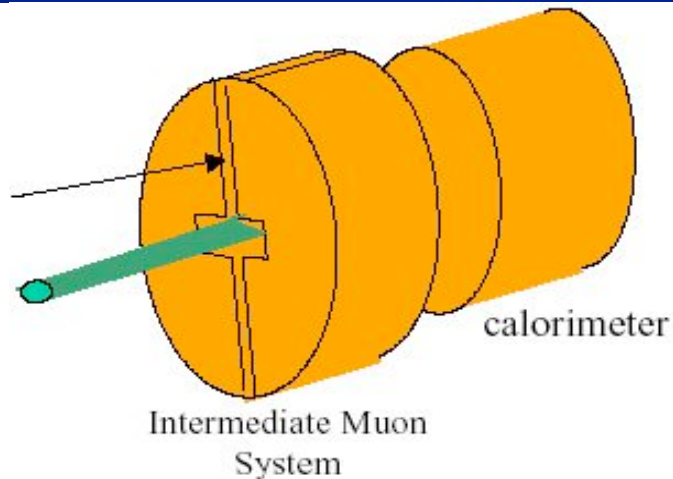
- Missing  $E_T$  distribution subject to many experimental effects
  - If anything goes wrong it will affect missing  $E_T$

# Sources of Instrumental Background

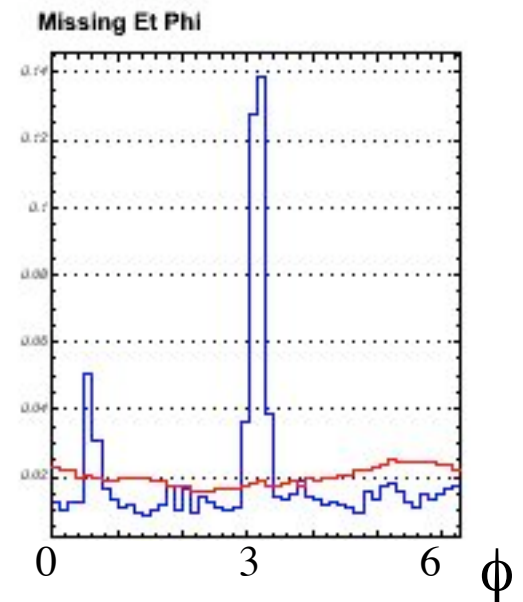
- Calorimeter Noise
  - Hot cells / coherent noise
    - Usually localized and can be rejected
- Calorimeter dead regions
  - Should only happen rarely in some runs
    - Should be removed by DQ criteria
- Cosmic rays and beam halo muons showering hard in calorimeter
  - Usually have no vertex but can overlap with MinBias event
    - Then have small tracking activity compared to calorimeter activity
  - Shower often only in hadronic calorimeter
- Example handles:
  - Track/calorimeter matches
  - Is direction of missing energy uniform?



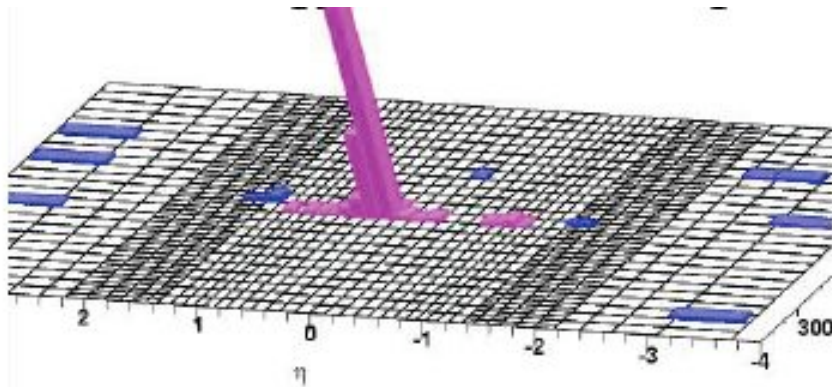
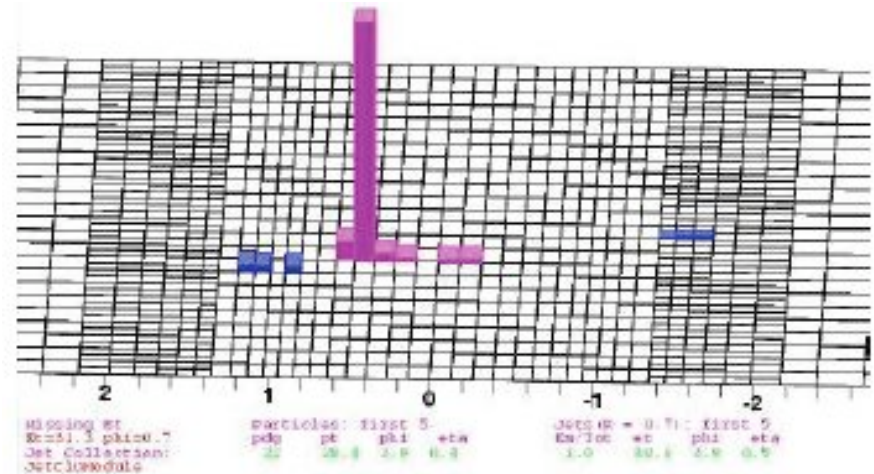
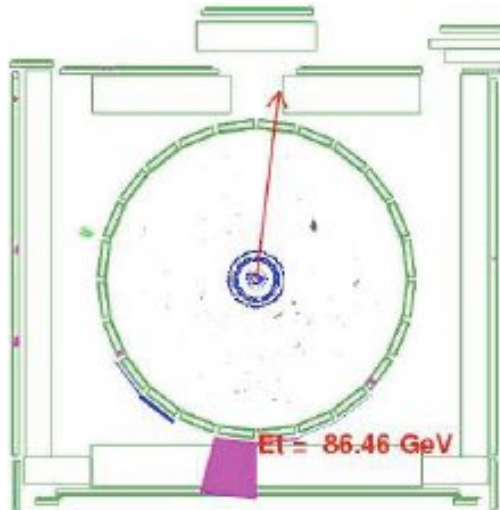
# Beam-Halo Muon Background



- Muon that comes from beam and goes through shielding
- Can cause showers in calorimeters
  - Shower usually looks not very much like physics jet
    - Often spike at certain azimuthal angles:  $\pi$
  - But there is lots of those muons!
  - Can even cause problem for trigger rate



# Some Cosmics and Beam-halo events



- Bigger problem for mono-jet than for multi-jet searches
- Can use
  - topological filters to reject events
  - Track matching calorimeter cluster

# Physics Backgrounds

- QCD multi-jet (mostly for 0-lepton case)
  - Missing  $E_T$  due to
    - Poor jet resolution / cracks in calorimeters
    - Neutrino momentum in semi-leptonic b/c- decays
- W/Z+jets
  - Missing  $E_T$  due to  $\nu$ 's from  $Z \rightarrow \nu\nu$ ,  $W \rightarrow l\nu$
- Top
  - Missing  $E_T$  due to  $\nu$ 's from  $t\bar{t} \rightarrow WbWb \rightarrow l\nu + X$

How do we estimate them?

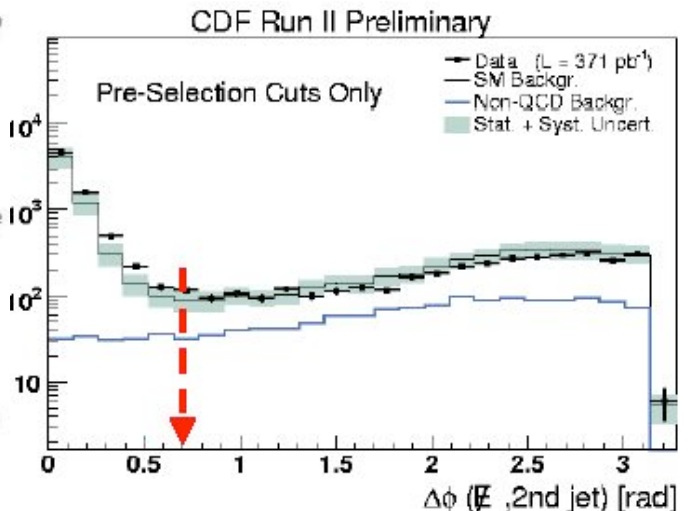
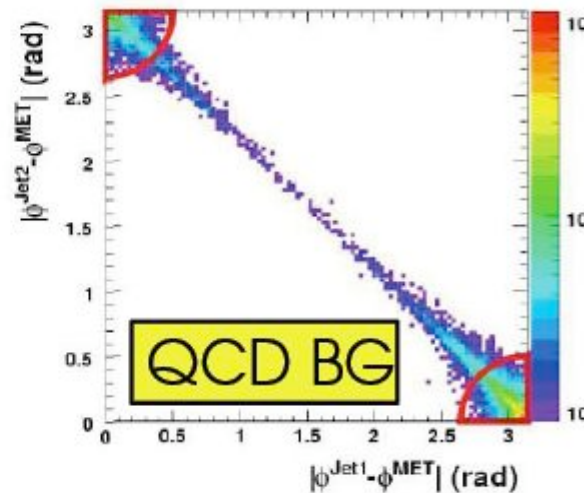
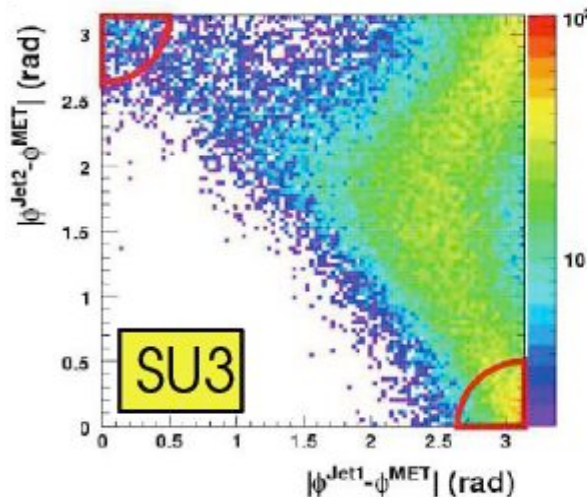
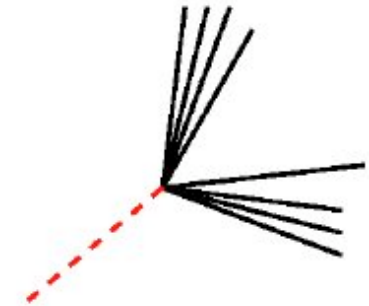
# QCD Multi-jet

- Require large  $\Delta\phi$ 
  - Between missing  $E_T$  and jets and between jets
  - Suppresses QCD dijet background due to jet mismeasurements

Background-like:  
 $\Delta\phi(\text{jet}, E_T^{\text{miss}}) \sim 0$



Signal-like:  
 $\Delta\phi(\text{jet}, E_T^{\text{miss}}) \gg 0$



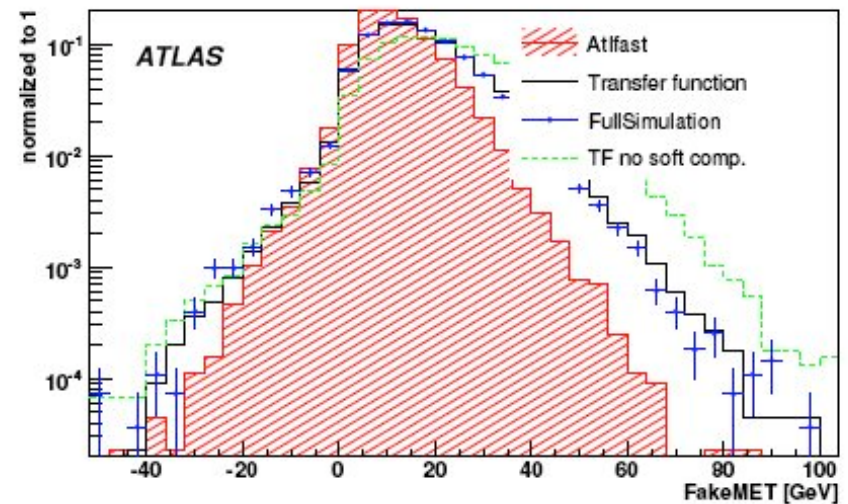
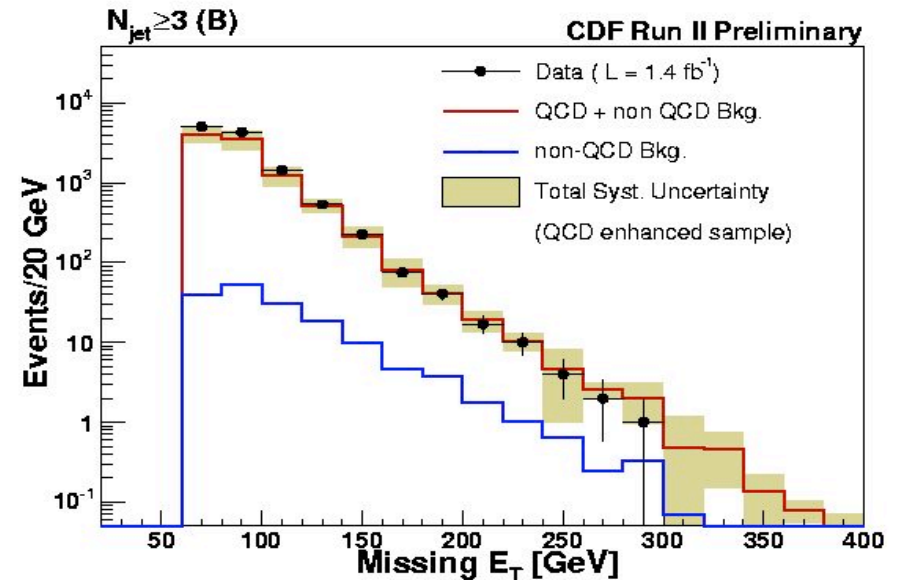
# Methods to estimate remaining QCD multi-jet Background

## 1. CDF uses MC

- Validate in region of low  $\Delta\Phi$  and low MET
- Extrapolate to large using MC
- Problem:
  - Relies on full MC simulation which can take “forever”

## 2. Parameterize truth jets with response function from full simulation

- Validate against full simulation
- Validate in region of lower MET
- Advantage:
  - Do not need to simulate as many events
- Need to make sure though that parameterization is really working

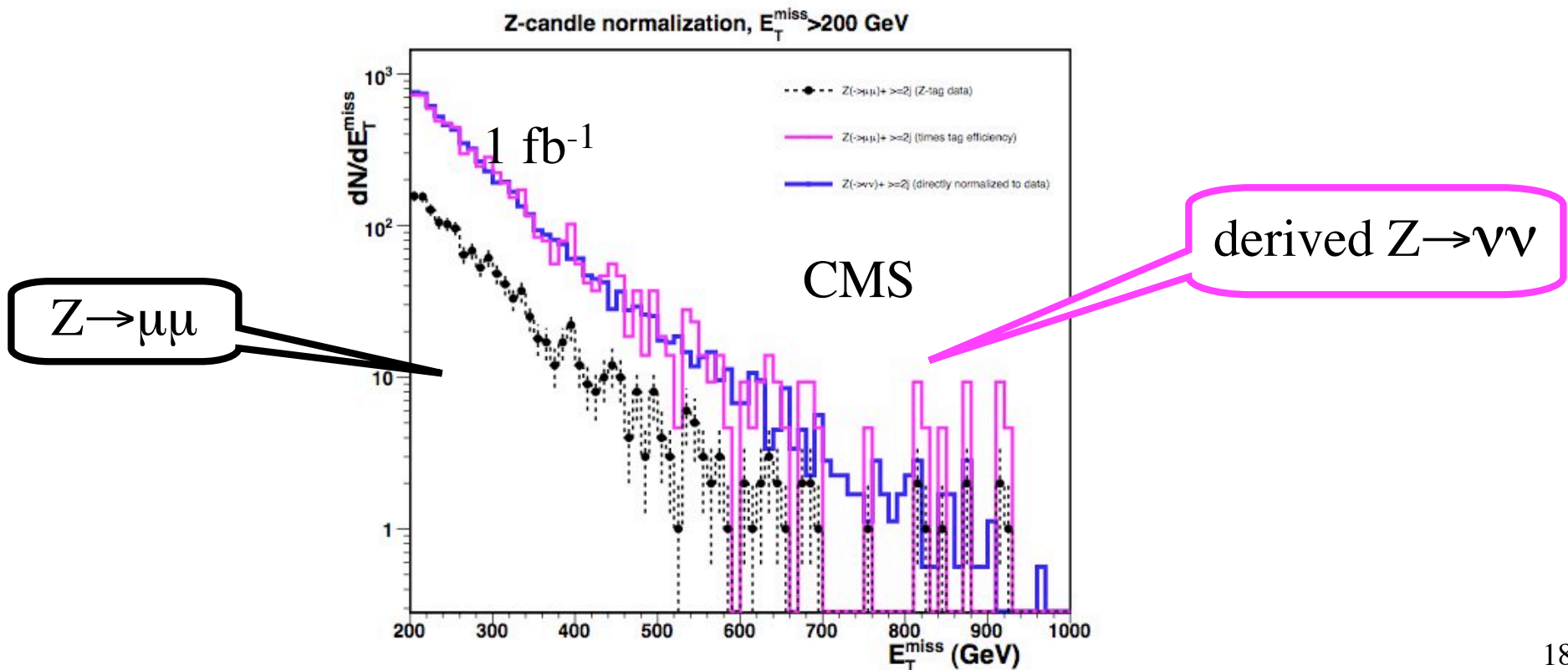




# Using $Z(\rightarrow \ell\ell)$ +jets for estimating $W/Z$ +jet background

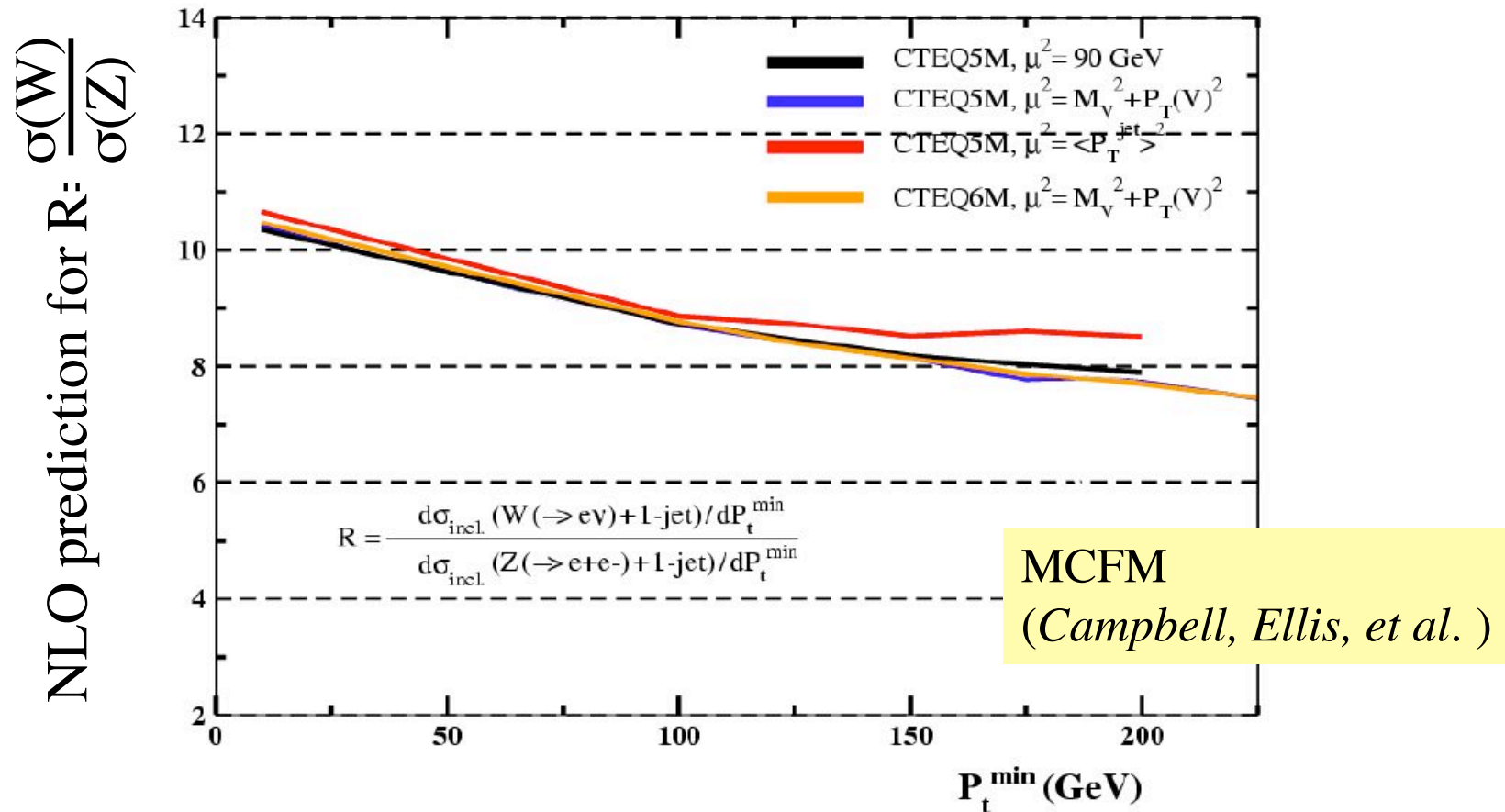
- Use  $Z(\rightarrow \ell\ell)$ +jets to extrapolate to  $Z(\rightarrow \nu\nu)$ +jets
  - $ME_T \sim p_T(Z)$

$$N_{Z \rightarrow \nu\bar{\nu}}(E_T^{\text{miss}}) = N_{Z \rightarrow \ell^+\ell^-}(p_T(\ell^+\ell^-)) \times c_{\text{Kin}}(p_T(Z)) \times c_{\text{Fidu}}(p_T(Z)) \times \frac{\text{Br}(Z \rightarrow \nu\bar{\nu})}{\text{Br}(Z \rightarrow \ell^+\ell^-)},$$



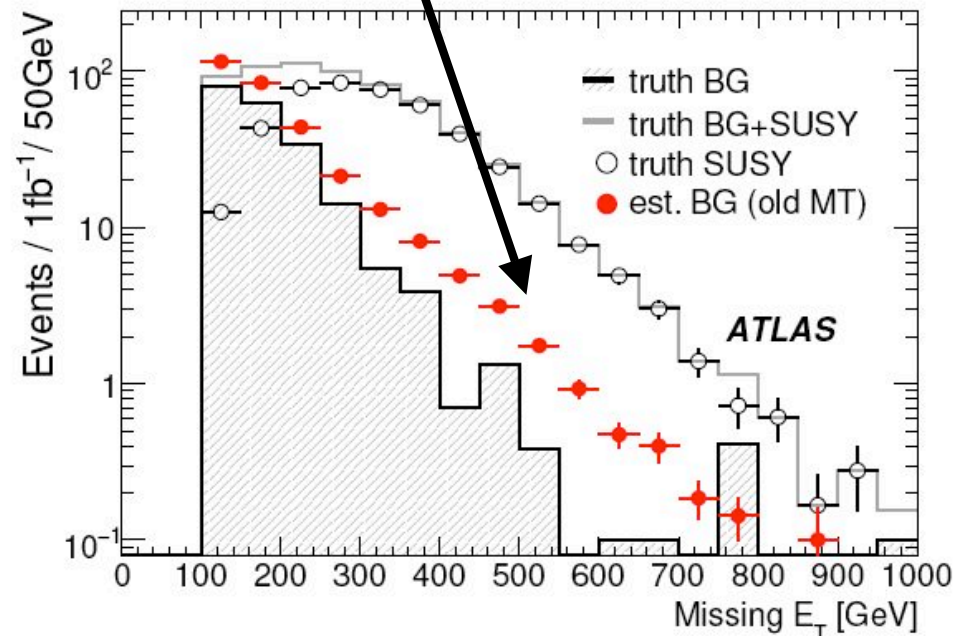
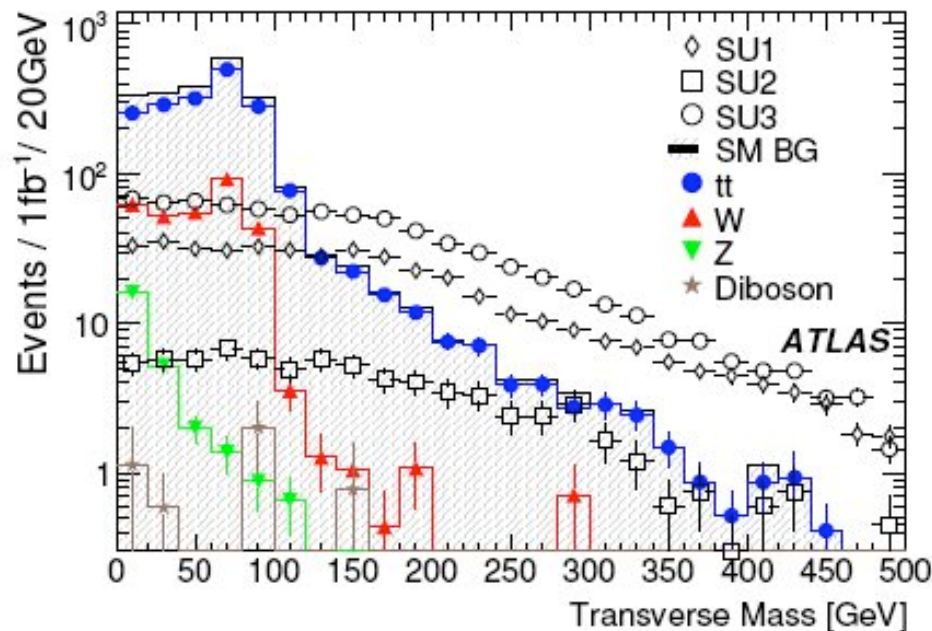
# W+jets background estimate

- Use Z→ll +jets also for this background too
  - Rely on theoretical prediction for W+jets vs Z+jets
    - This is well known though (<15%)!



# Top and W+jets background estimate

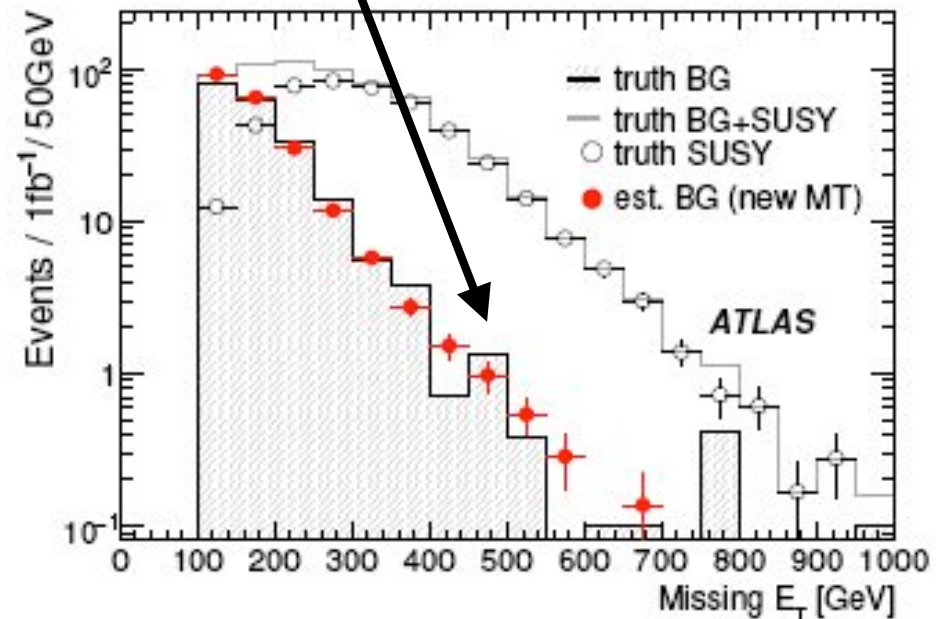
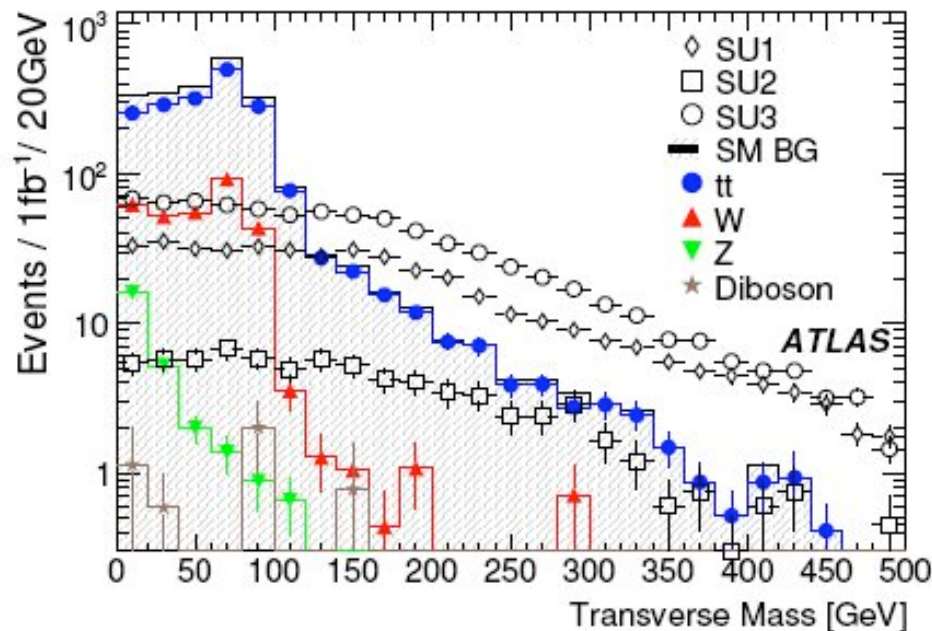
- Use region of low  $m_T(W)$ 
  - Extrapolate to signal region using MC
  - But may be contaminated by SUSY => overestimate BG
    - depending on specifics of model





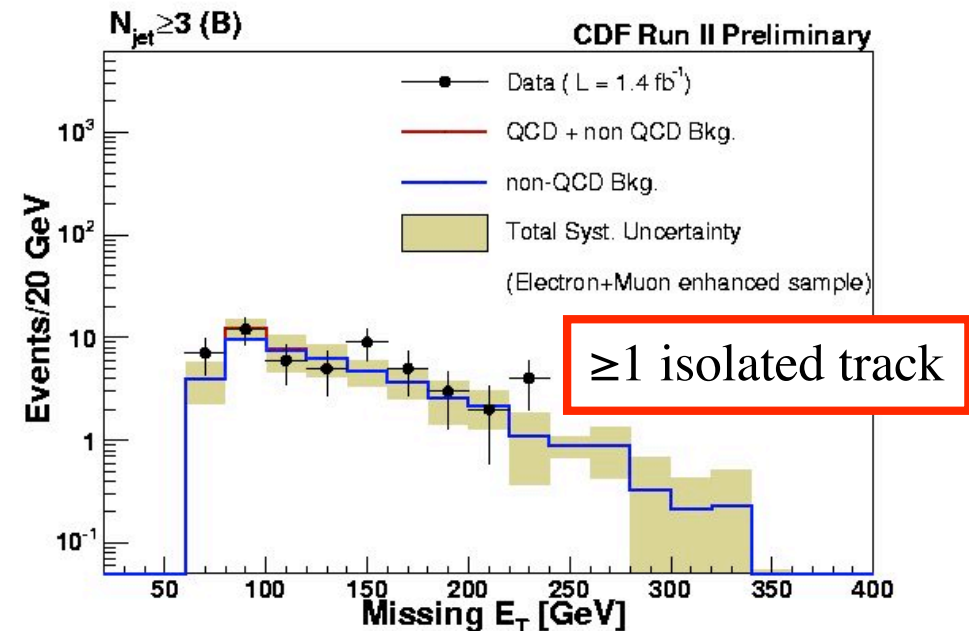
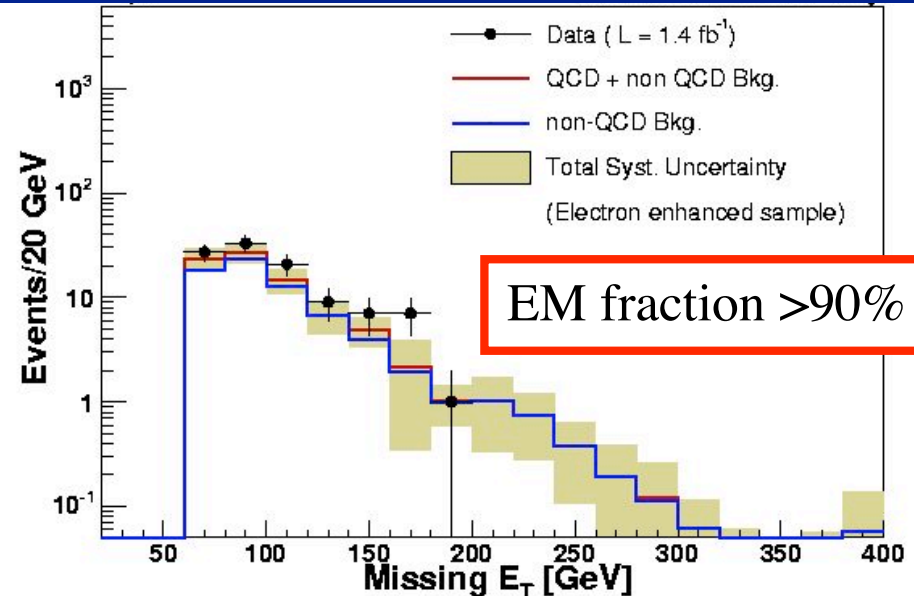
# Top and W+jets background estimate

- Use region of low  $m_T(W)$ 
  - Extrapolate to signal region using MC
  - But may be contaminated by SUSY => overestimate
    - depending on specifics of model
  - Can attempt “SUSY background subtraction” to correct for it

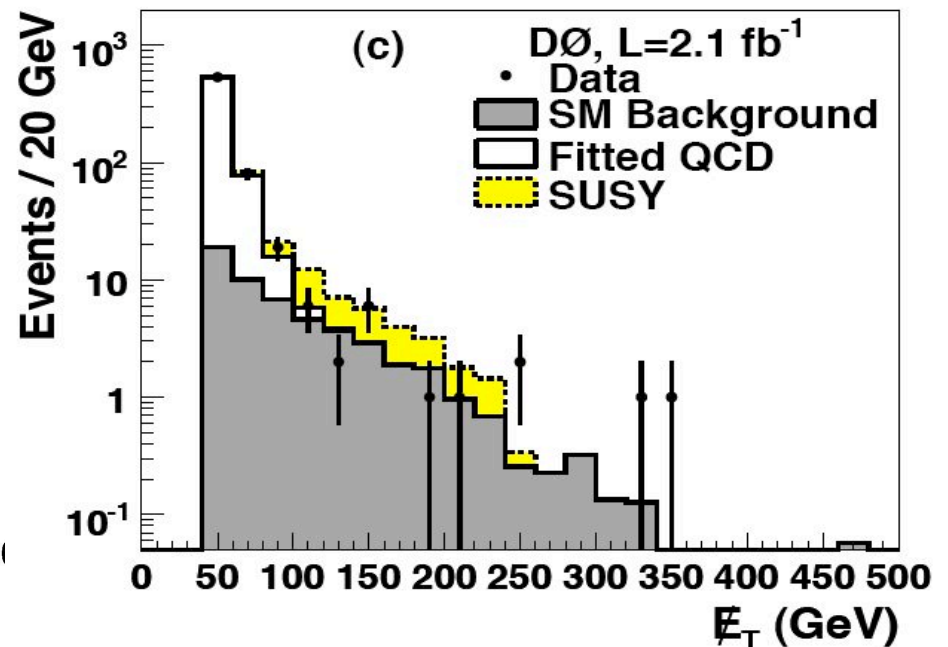
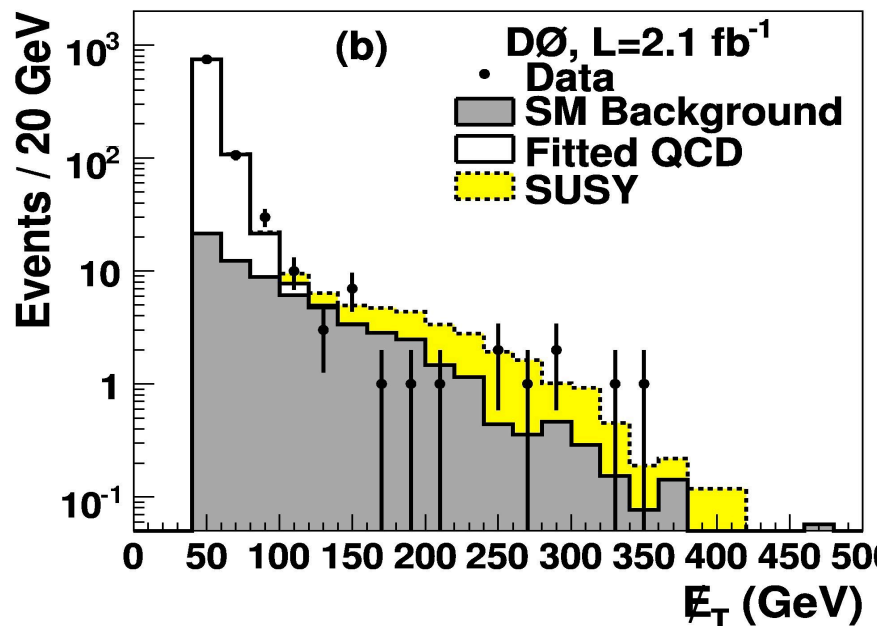
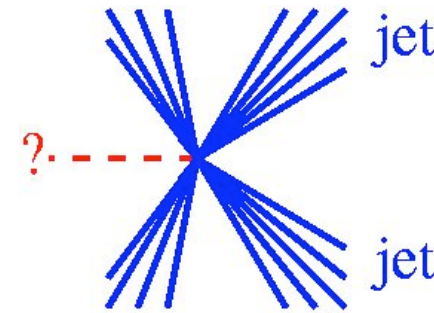
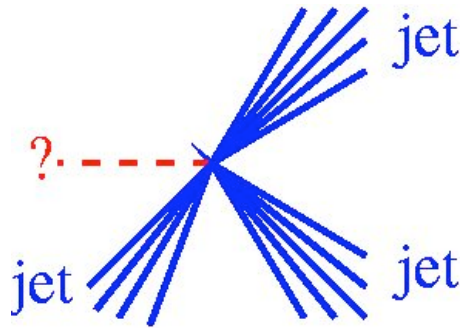


# W+jets, Z+jets and Top background

- Checks at Tevatron 0-lepton analysis
  - Background sources:
    - W/Z+jets, top
    - Suppressed by vetoes:
      - Events with jet with EM fraction > 90%
        - » Rejects electrons
      - Events with isolated track
        - » Rejects muons, taus and electrons
  - Define control regions:
    - W/Z+jets, top
      - Make all selection cuts but invert lepton vetoes
    - Gives confidence in those background estimates
      - Modeled using Alpgen MC
      - Cross sections determined using NLO calculation
- May not work at LHC due to expectation of large cascade decays

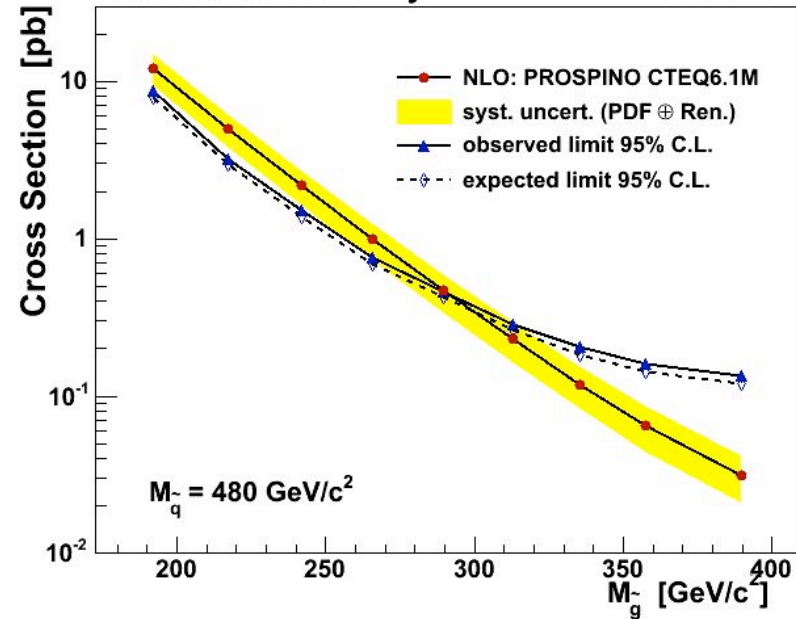
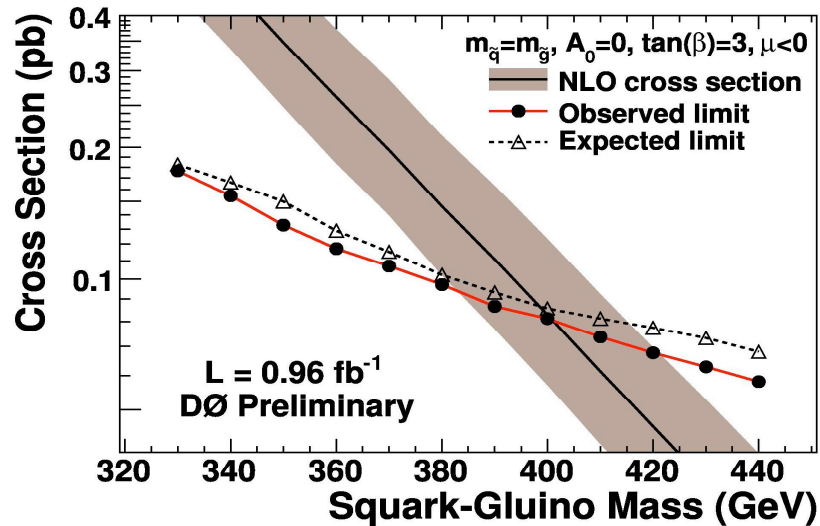


# Final Analysis Plots at the Tevatron



Data agree with background estimate => derive limits

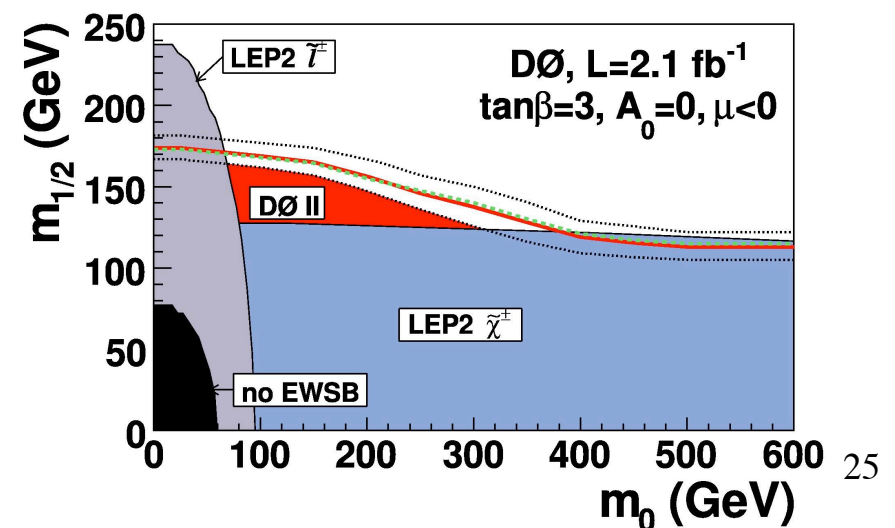
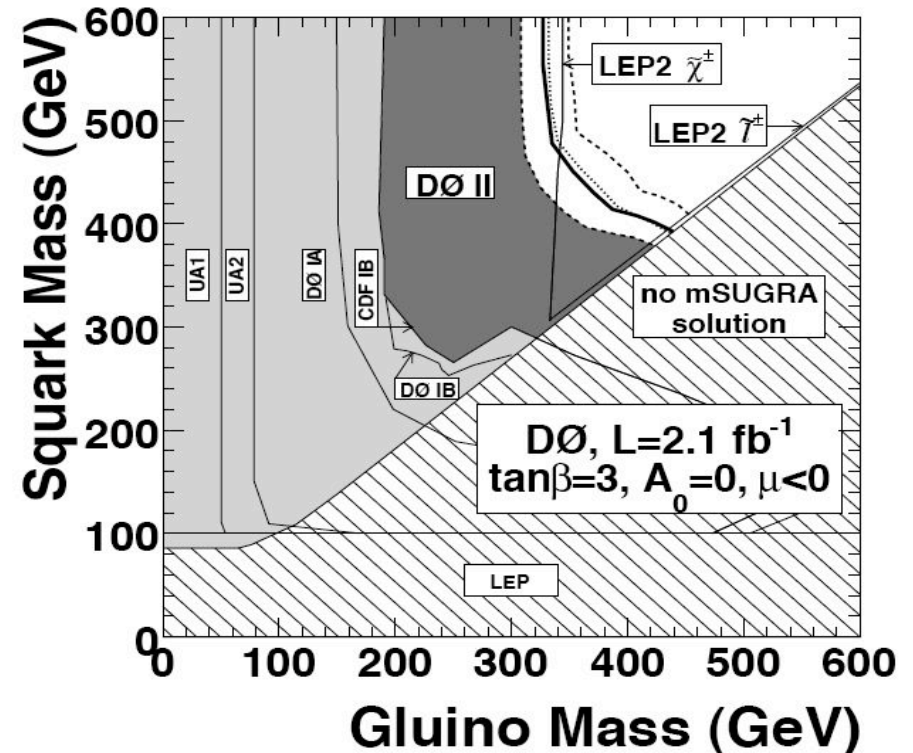
# Cross Section Limits



- No excess in data
  - Evaluate upper limit on cross section
  - Find out where it crosses with theory
- Theory has large uncertainty:  $\sim 30\%$ 
  - Crossing point with theory lower bound  $\sim$  represents limit on squark/gluino mass

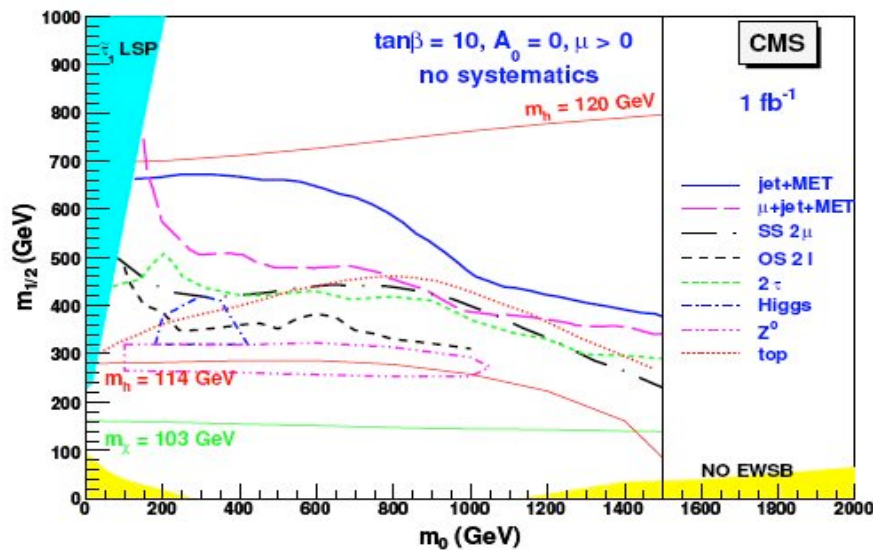
# Squark and Gluino Mass Limits

- Set constraints on masses at EWK scale:
  - $M(\tilde{g}) > 308 \text{ GeV}$
  - $M(\tilde{q}) > 379 \text{ GeV}$
- Can also be represented in terms of GUT scale parameters
  - Within constrained models

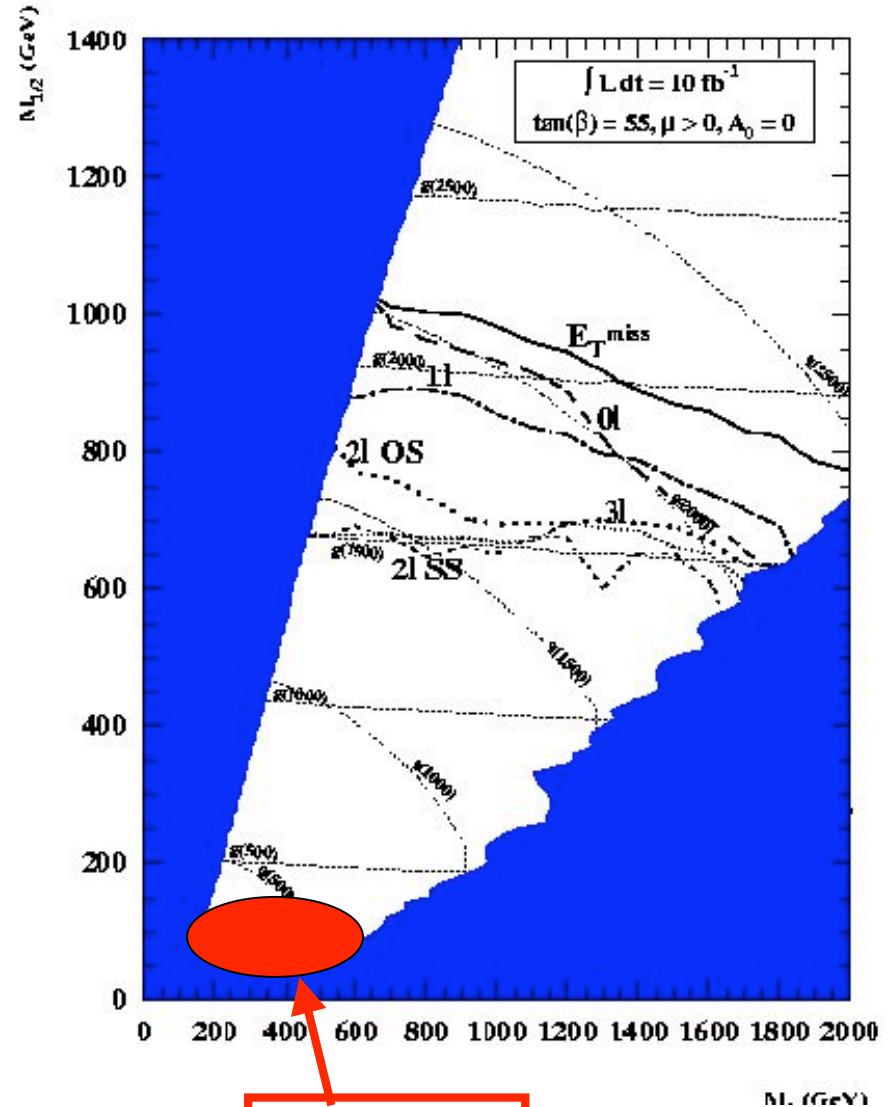




# LHC SUSY Discovery Reach



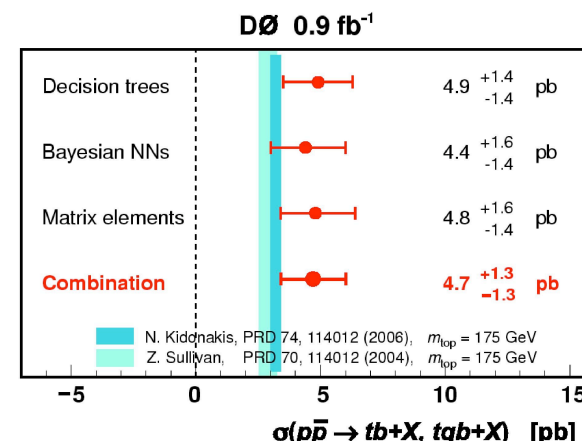
- With  $1 \text{ fb}^{-1}$ :
  - Sensitive to  $m(\tilde{g}) < 1000 \text{ GeV}/c^2$
- With  $10 \text{ fb}^{-1}$ :
  - Sensitive to  $m(\tilde{g}) < 1800 \text{ GeV}/c^2$
- Amazing potential!
  - If data can be understood
  - If current MC predictions are  $\approx$ ok



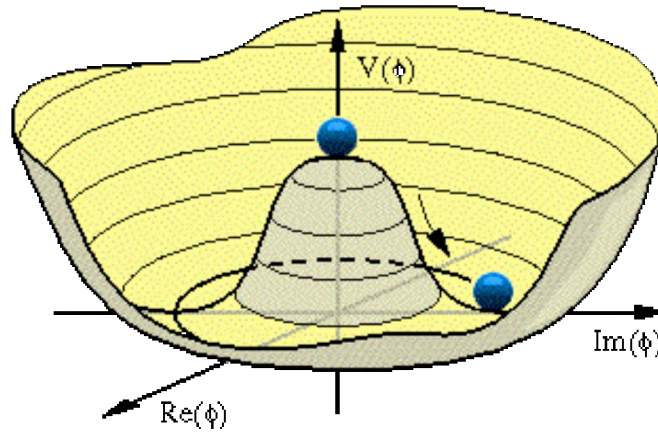
Tevatron

# Some Remarks on Advanced Analysis Techniques

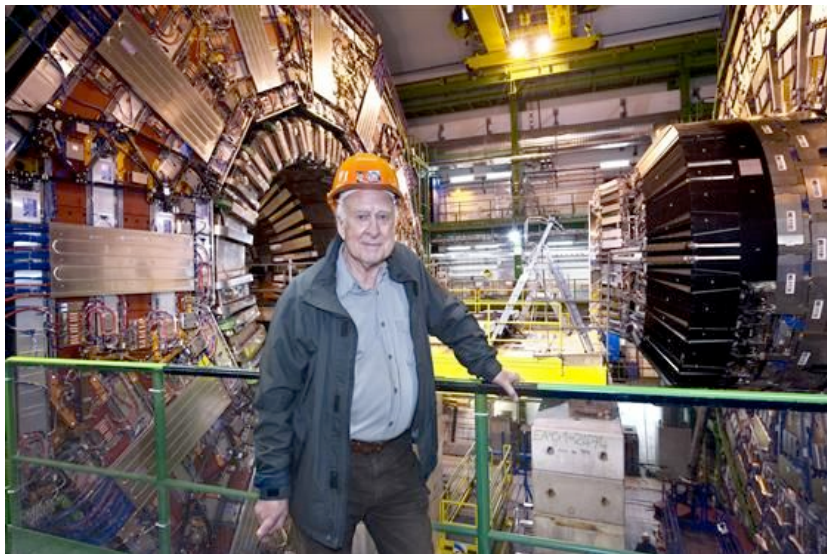
- **Quite a few techniques available:**
  - Neural Network, Likelihood, Boosted Decision Tree, Matrix Element, ...
  - No clear winner has yet been identified
    - Generally NN least transparent probably



- **Why do we trust them less than simple analyses?**
  - Simple kinematic quantities can be calculated at NLO by theorists while e.g. NN distribution cannot
    - Gives confidence, good cross-check!
  - Techniques exploit correlations between variables
    - Harder to understand if the MC models correlations correctly
    - More validation needed ( $\Rightarrow$  analysis takes longer)
  - Less transparent
    - Worry is always that it exploits some MC feature that does not reflect the data
- **Can and has been done of course though**
  - But only in mature experiments



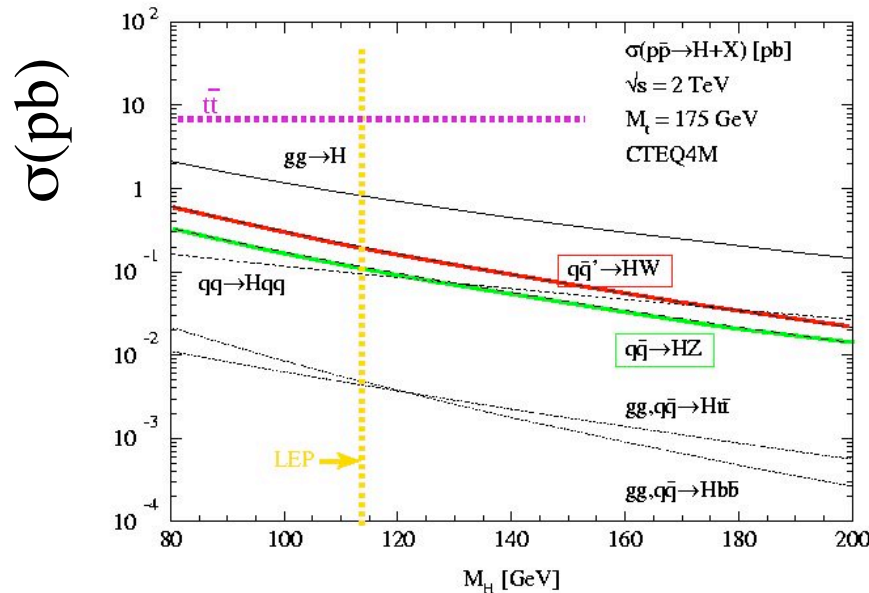
# The Higgs Boson



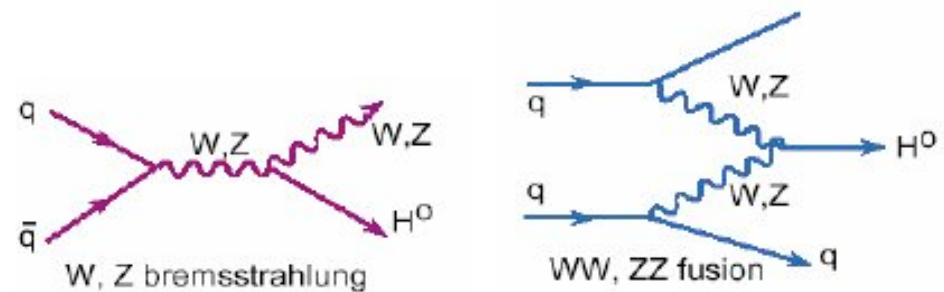
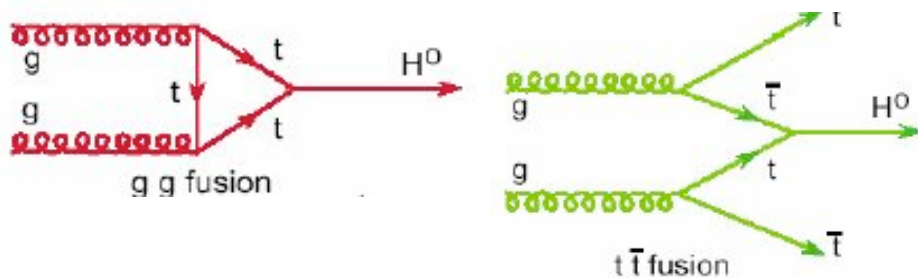
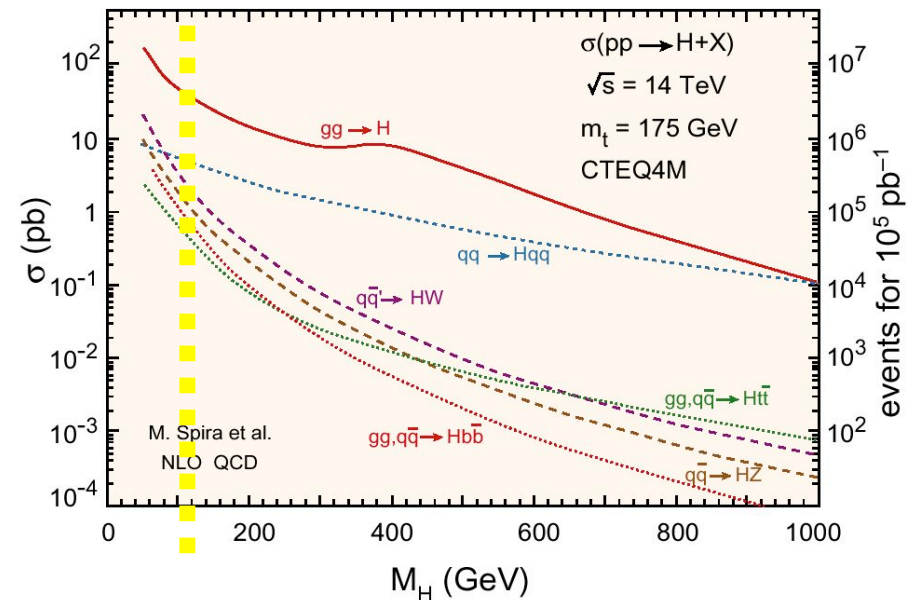


# Higgs Production: Tevatron and LHC

Tevatron



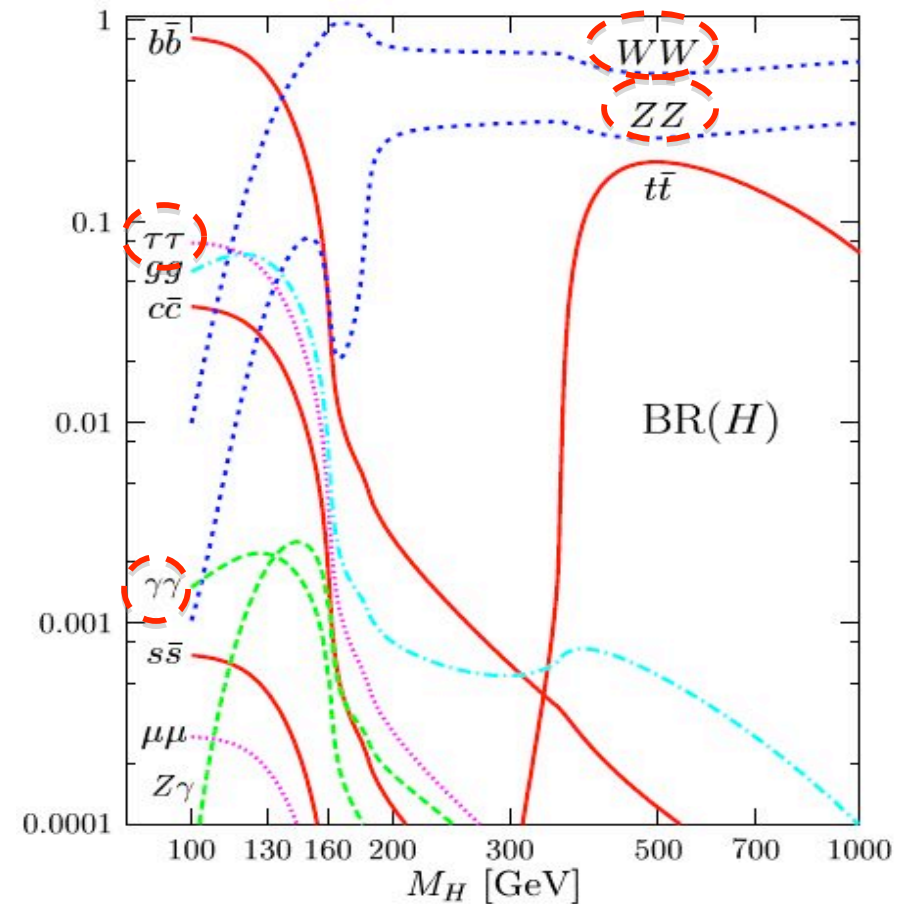
LHC



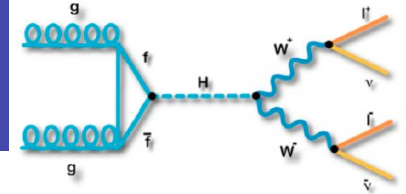
dominant:  $gg \rightarrow H$ , subdominant:  $HW$ ,  $HZ$ ,  $Hq\bar{q}$

# Higgs Boson Decay

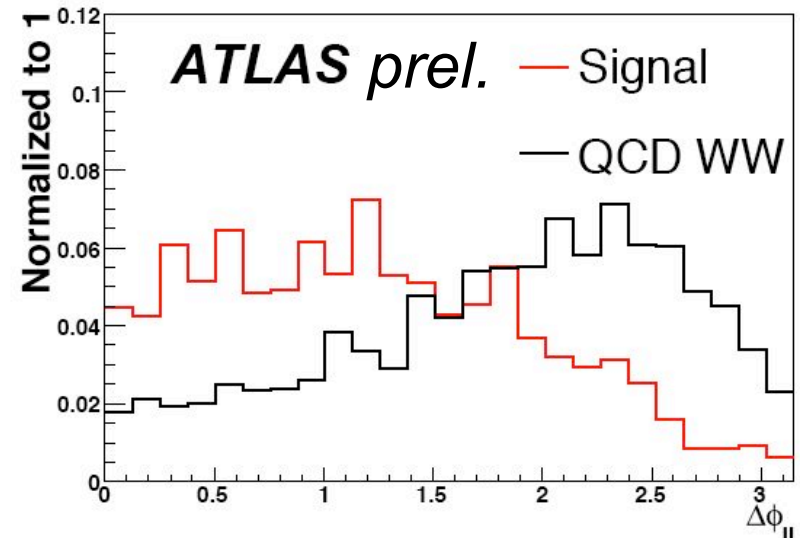
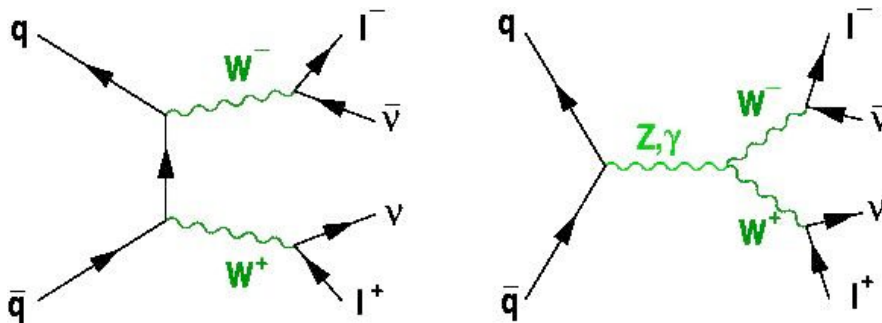
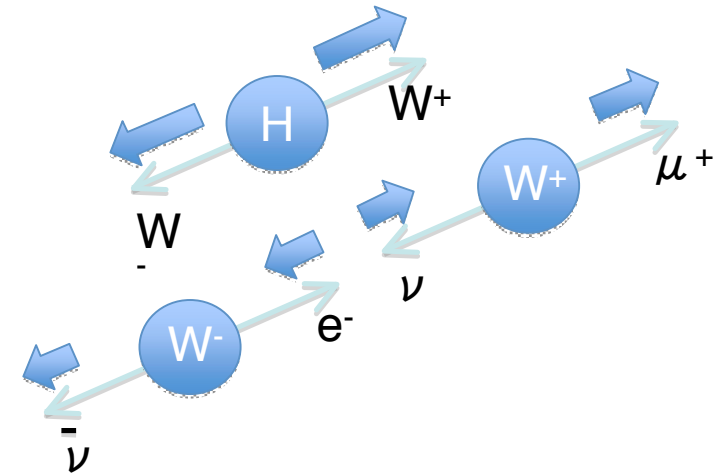
- Depends on Mass
- $M_H < 130 \text{ GeV}/c^2$ :
  - $b\bar{b}$  dominant
  - $WW$  and  $\tau\tau$  subdominant
  - $\gamma\gamma$  small but useful
- $M_H > 130 \text{ GeV}/c^2$ :
  - $WW$  dominant
  - $ZZ$  cleanest



$$\mathbf{H} \rightarrow \mathbf{WW}^{(*)} \rightarrow \mathbf{l}^+ \mathbf{l}^- \bar{\nu} \nu$$



- Higgs mass reconstruction impossible due to two neutrinos in final state
- Make use of spin correlations to suppress WW background:
  - Higgs is scalar: spin=0
  - leptons in  $H \rightarrow WW^{(*)} \rightarrow l^+ l^- \nu \bar{\nu}$  are collinear
- Main background:
  - WW production



# $H \rightarrow WW^{(*)} \rightarrow l^+ l^- \nu \nu \quad (l=e, \mu)$

- **Event selection:**

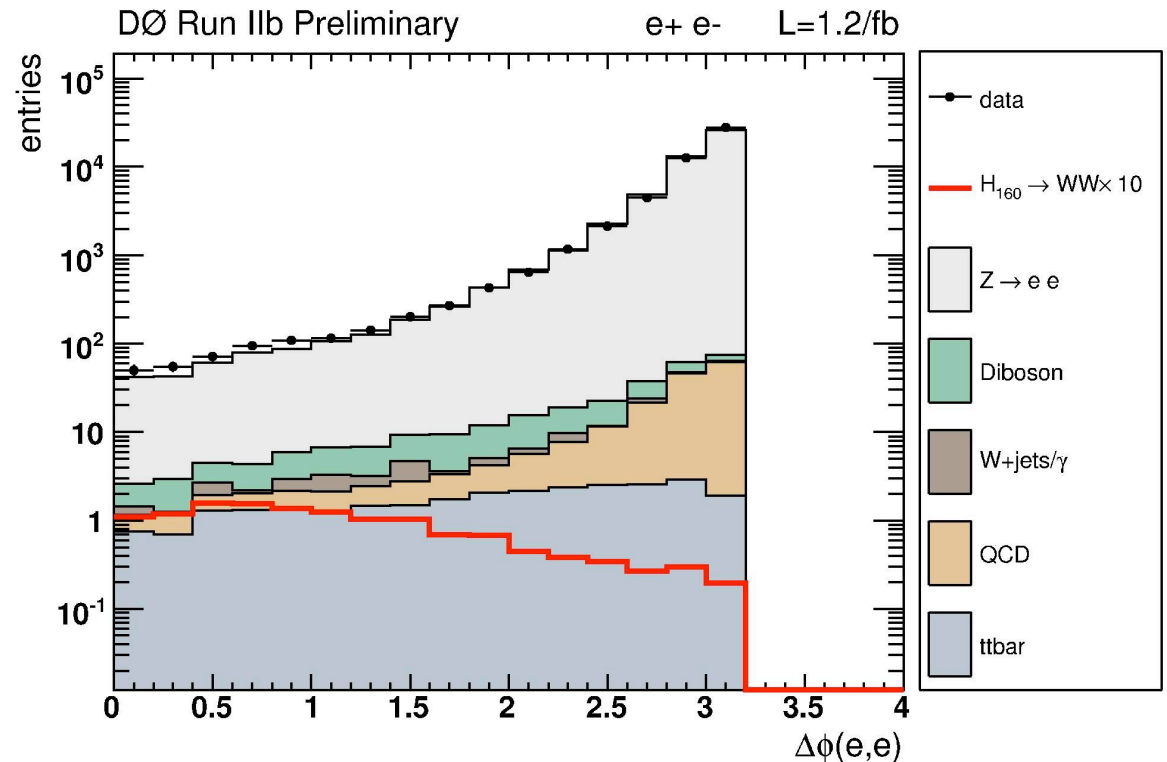
- 2 isolated  $e/\mu$  :
  - $p_T > 15, 10$  GeV
- Missing  $E_T > 20$  GeV
- Veto on
  - Z resonance
  - Energetic jets

- **Main backgrounds**

- SM  $WW$  production
- Top
- Drell-Yan
- Fake leptons

- **Plot everything under the sun**

- to convince yourself you have the background right

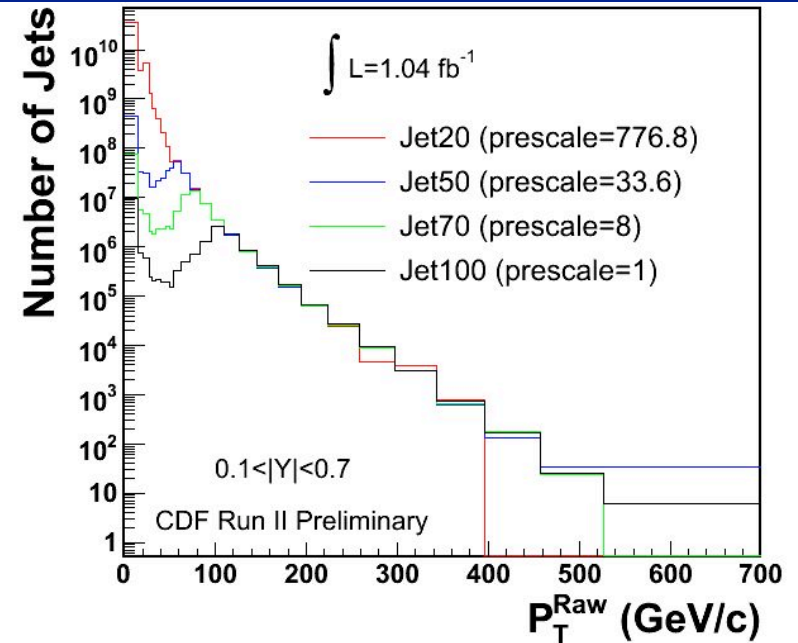


# Jets faking Electrons

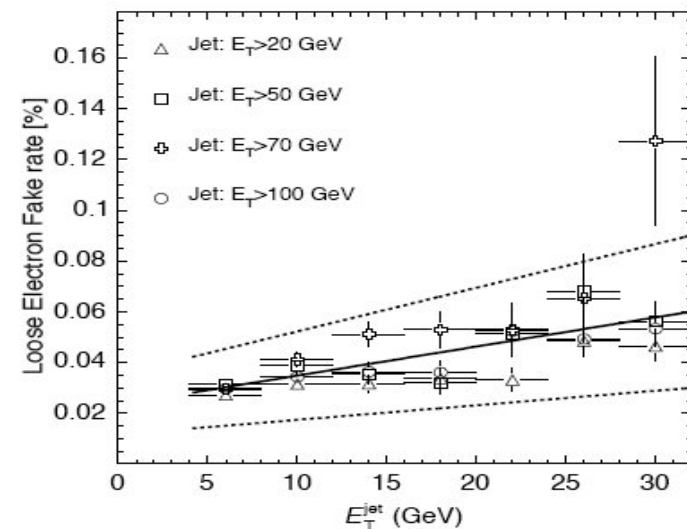
- Jets can pass electron ID cuts,
  - Mostly due to
    - early showering charged pions
    - Conversions:  $\pi^0 \rightarrow \gamma\gamma \rightarrow ee + X$
    - Semileptonic b-decays
  - ♣ Difficult to model in MC
    - Hard fragmentation
    - Detailed simulation of calorimeter and tracking volume
- Measured in inclusive jet data at various  $E_T$  thresholds
  - Prompt electron content negligible:
    - $N_{\text{jet}} \sim 10$  billion at 50 GeV!
  - Fake rate per jet:**

	CDF	ATLAS
Loose cuts	$5 \times 10^{-4}$	$5 \times 10^{-3}$
Tight cuts	$1 \times 10^{-4}$	$1 \times 10^{-5}$

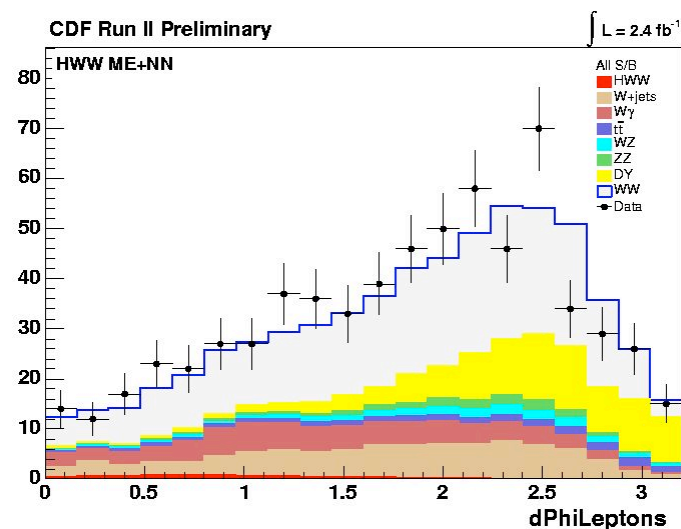
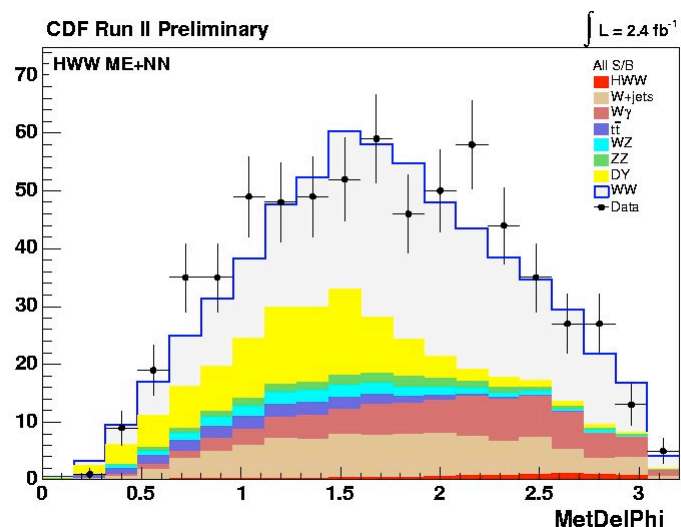
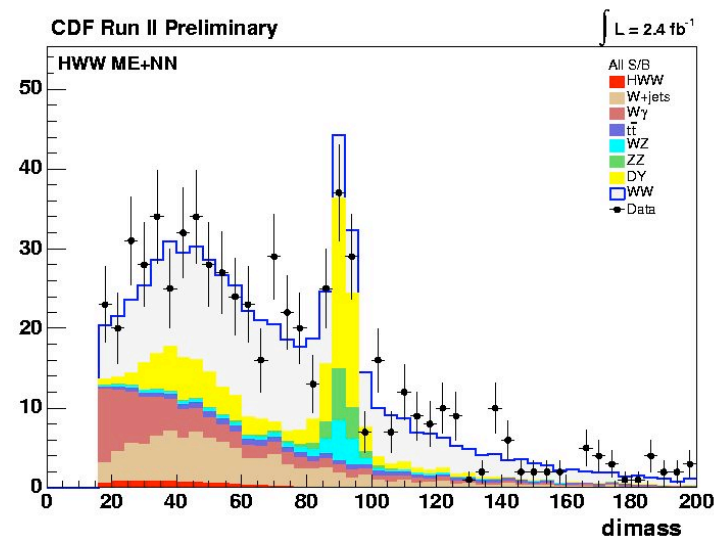
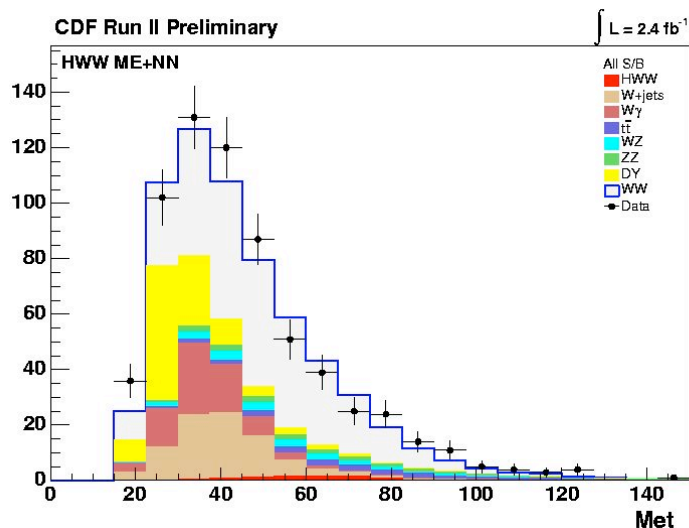
- Typical uncertainties 50%



Fake Rate (%)



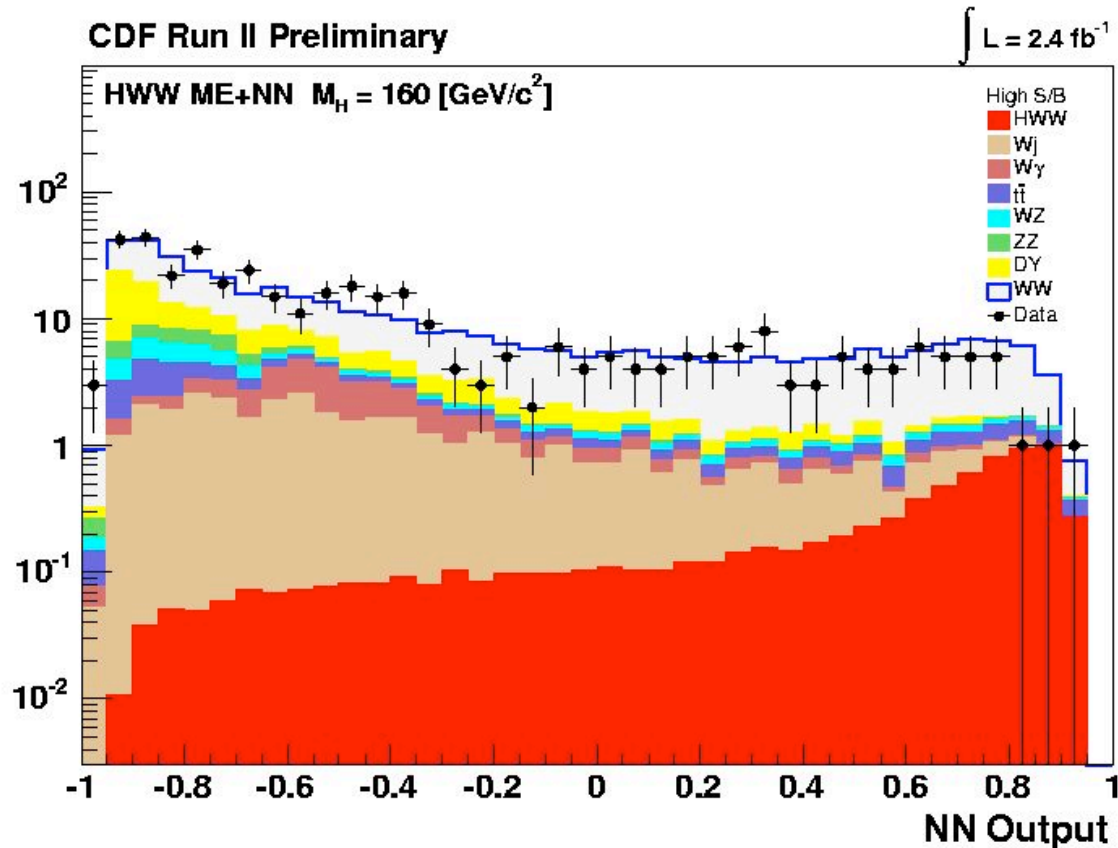
# Plot Everything Under the Sun..



- Validates the background prediction
  - Very often these plots “don’t work” since there is some problem
  - Now plug all into sophisticated techniques!

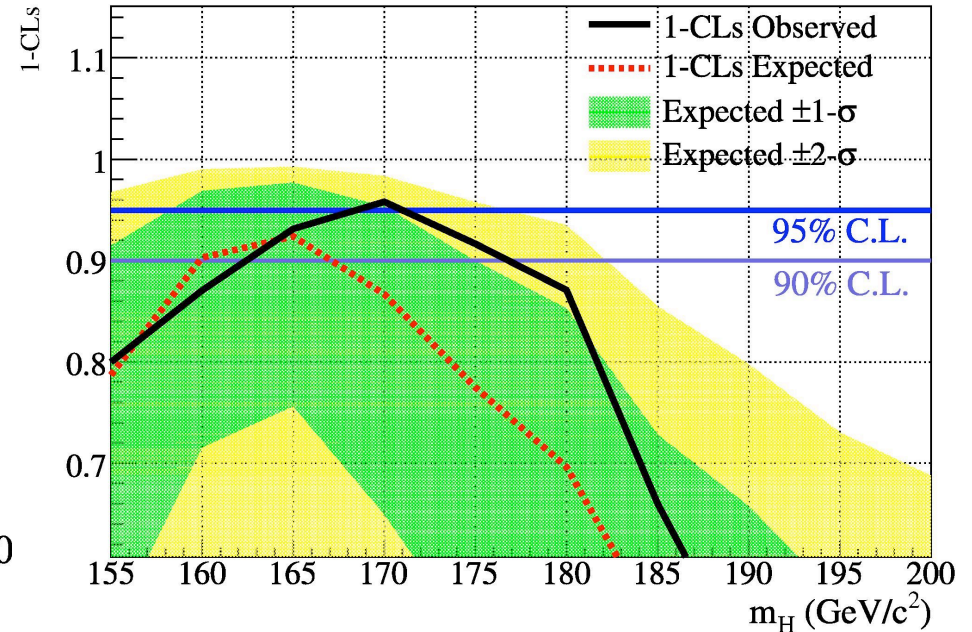
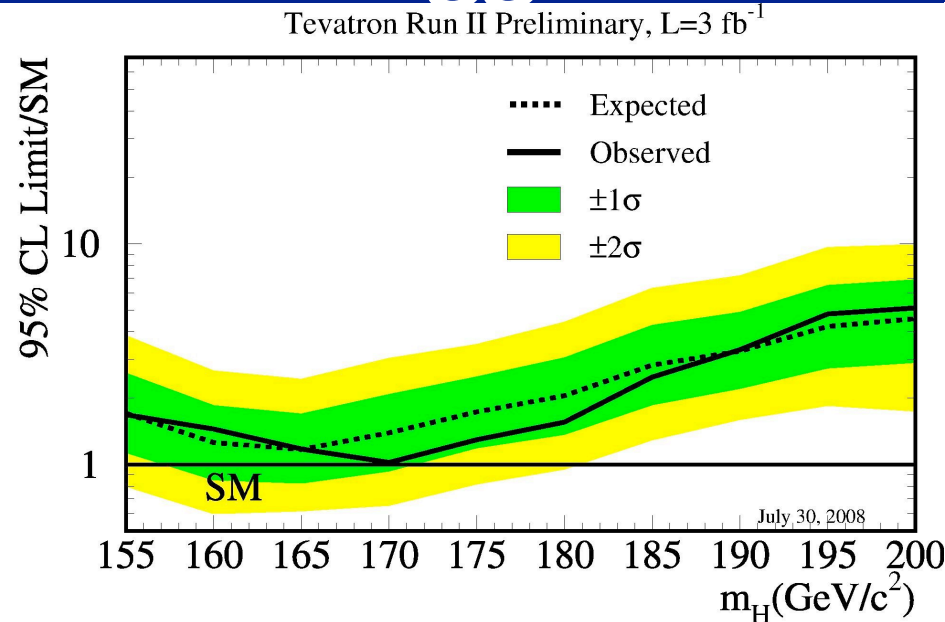


# Result of Neural Network



- At high NN Output data fluctuated low
  - No sign of Higgs boson

# Tevatron limits on the Higgs boson cross section

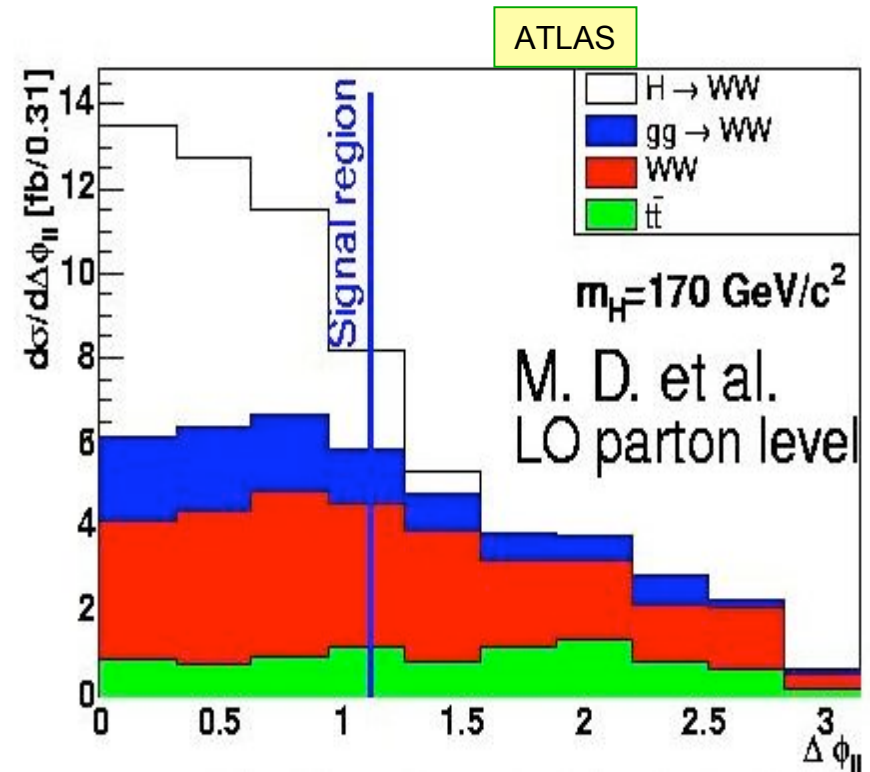
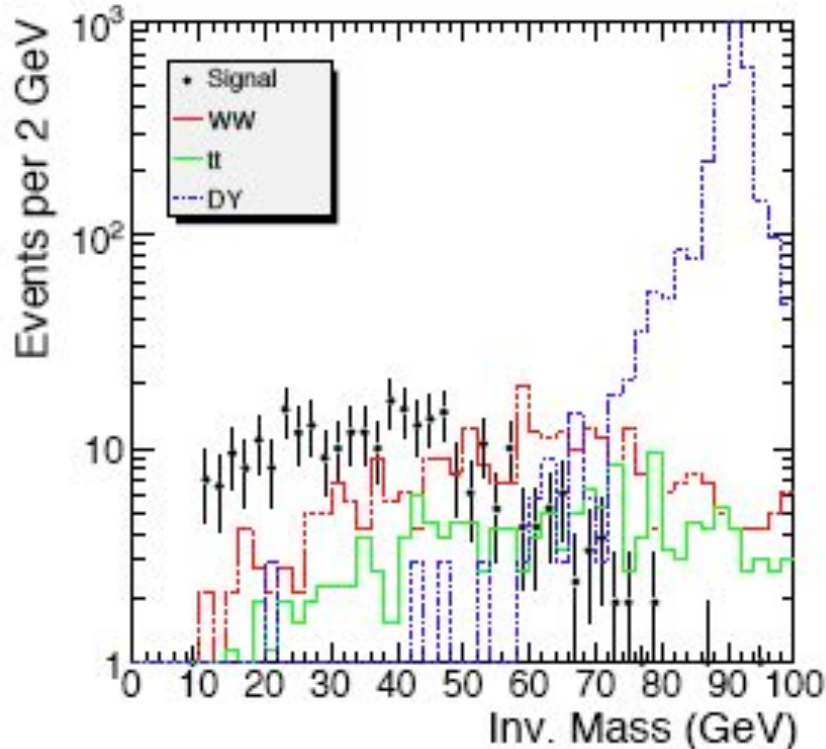


- Lack of observation
  - $\Rightarrow$  an upper limit on the Higgs cross section
  - I.e. if the cross section was large we would have seen it!
- Results presented typically as ratio:
  - Experimental limit / theoretical cross section
  - If this hits 1 we exclude the Higgs boson at that mass!
- Brand new result: exclude  $M_H=170 \text{ GeV}/c^2$



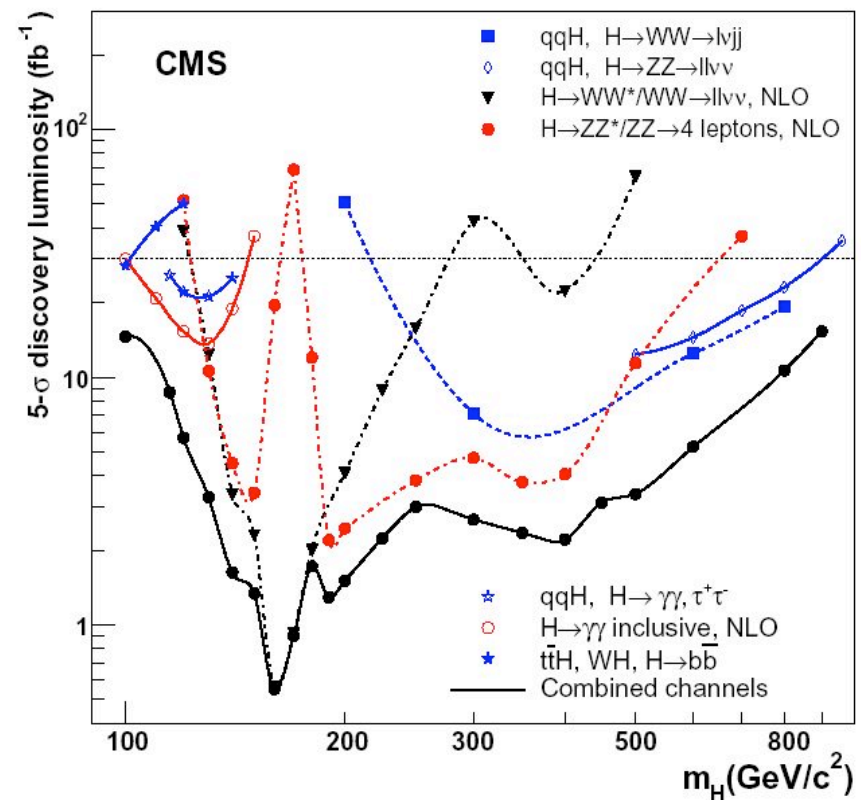
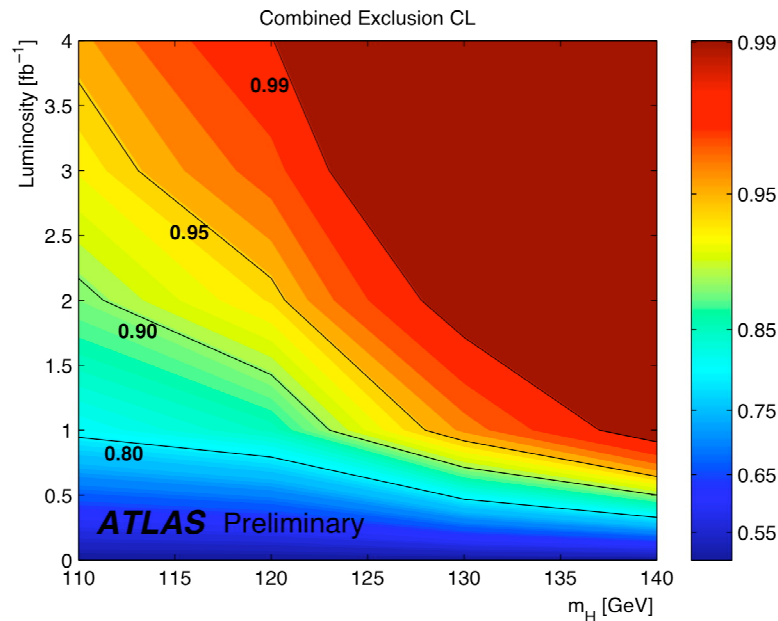
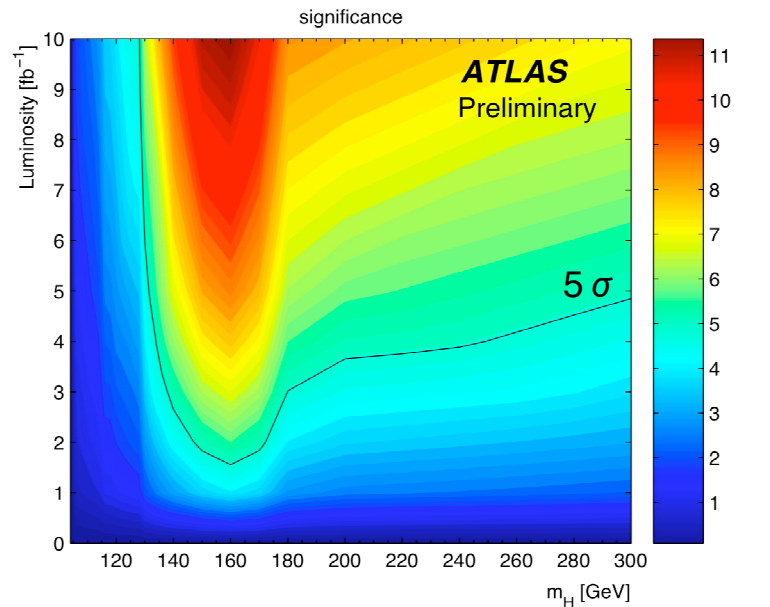
# Early Higgs Signals at LHC

$H \rightarrow WW^* (m_H=170 \text{ GeV})$



LHC has about 4 times better  
signal / background than Tevatron

# LHC SM Higgs Discovery Potential



- 5 $\sigma$  discovery over full mass range with  $\sim 20 \text{ fb}^{-1}$ 
  - Most challenging at low mass
- 95% exclusion over full mass range with  $\sim 4 \text{ fb}^{-1}$

# Concluding Remarks

## Data are very precious

- Treat them with the highest respect
- Try to not jump to conclusions too fast
  - Data analysis is like detective work
- Try to use all you can to understand them
  - Redundancy of detector
    - tracker vs calorimeter etc.
  - Complementary physics processes
    - W's vs Z's etc.
  - Monte Carlo tools and theoretical calculations
- Use your brain and your judgment

**This was my very personal view on the key issues  
of data analysis**

**Thanks,**  
**and lot's of fun and luck**  
**for your analyses**

# High Mass Resonances



# Resonances or Tails

- New resonant structure:

- New gauge boson:

- $Z' \rightarrow ee, \mu\mu, \tau\tau, tt$
    - $W' \rightarrow e\nu, \mu\nu, \tau\nu, tb$

- Randall-Sundrum Graviton:

- $G \rightarrow ee, \mu\mu, \tau\tau, \gamma\gamma, WW, ZZ, \dots$

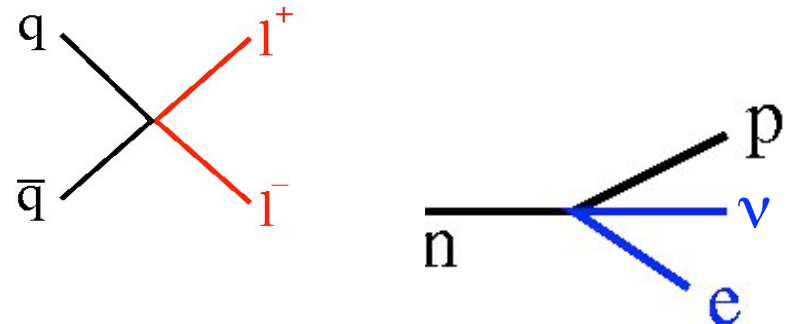
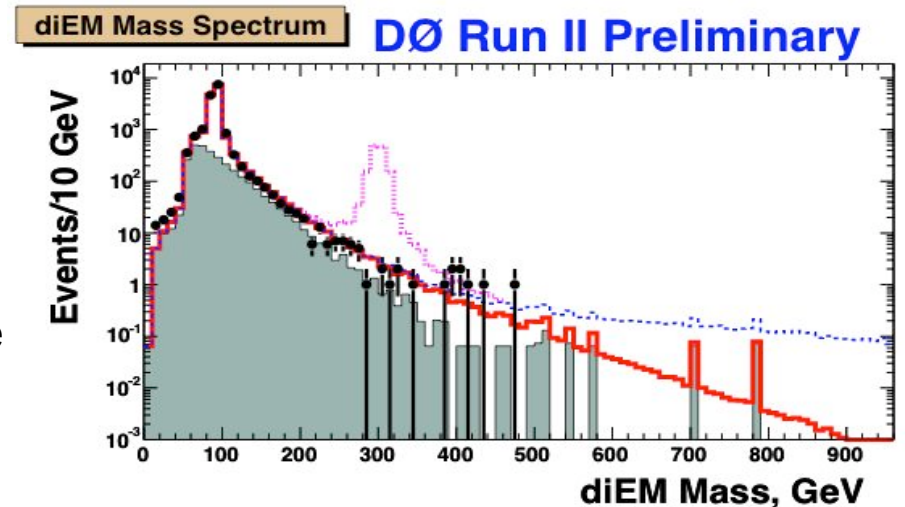
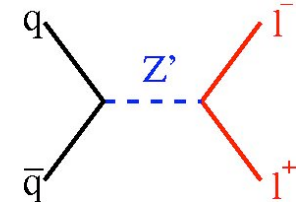
- Tail:

- Large extra dimensions (ADD model)

- Many many many resonances close to each other:
    - “Kaluza-Klein-Tower”:  $ee, \mu\mu, \tau\tau, \gamma\gamma, WW, ZZ, \dots$

- Contact interaction

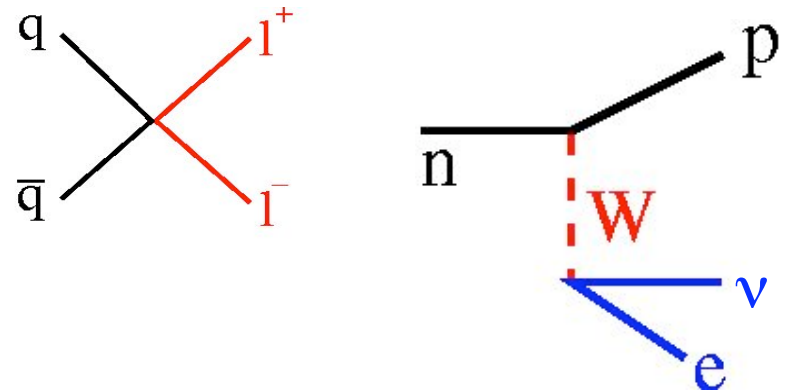
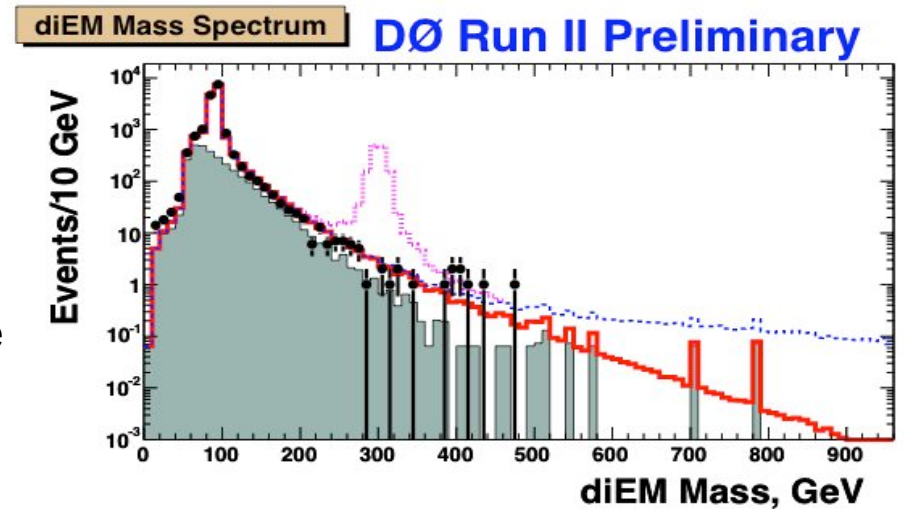
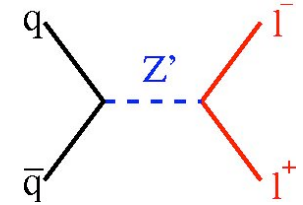
- Effective 4-point vertex
      - E.g. via t-channel exchange of very heavy particle
    - Like Fermi's  $\beta$ -decay





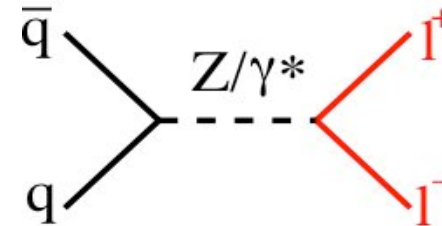
# Resonances or Tails

- New resonant structure:
  - New gauge boson:
    - $Z' \rightarrow ee, \mu\mu, \tau\tau, tt$
    - $W' \rightarrow e\nu, \mu\nu, \tau\nu, tb$
  - Randall-Sundrum Graviton:
    - $G \rightarrow ee, \mu\mu, \tau\tau, \gamma\gamma, WW, ZZ, \dots$
- Tail:
  - Large extra dimensions (ADD model)
    - Many many many resonances close to each other:
    - “Kaluza-Klein-Tower”:  $ee, \mu\mu, \tau\tau, \gamma\gamma, WW, ZZ, \dots$
  - Contact interaction
    - Effective 4-point vertex
      - E.g. via t-channel exchange of very heavy particle
    - Like Fermi’s  $\beta$ -decay

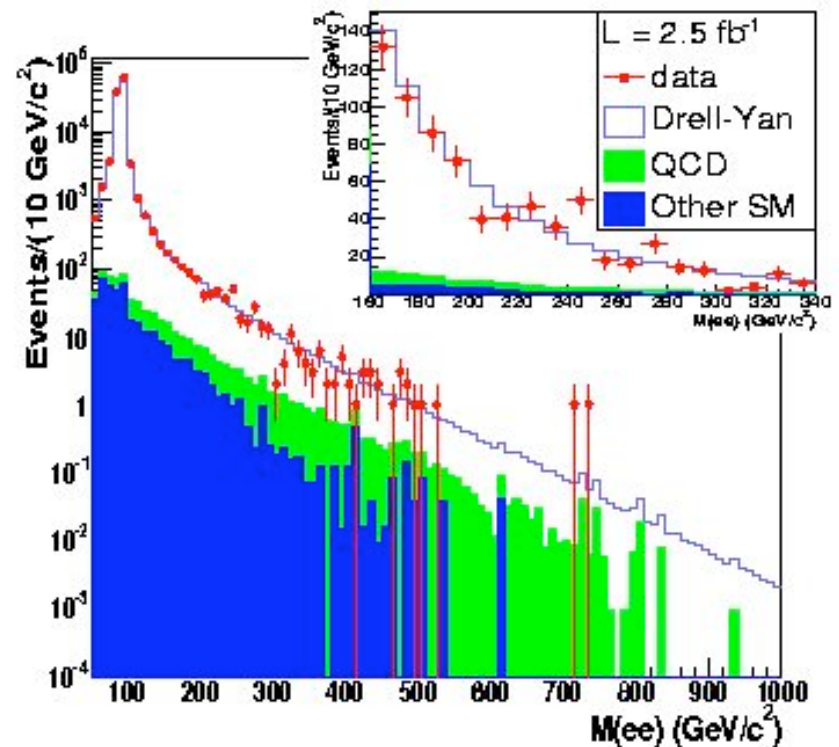


# Dilepton Selection

- Background to two high momentum leptons
  - irreducible background is Drell-Yan production
    - Well known/calculated by MC
    - NNLO calculation also available
  - Other backgrounds:
    - Jets faking leptons: reject by making optimal lepton ID cuts
    - WW, diphoton, etc. very small
- None of these backgrounds are expected to make peaks in mass distribution

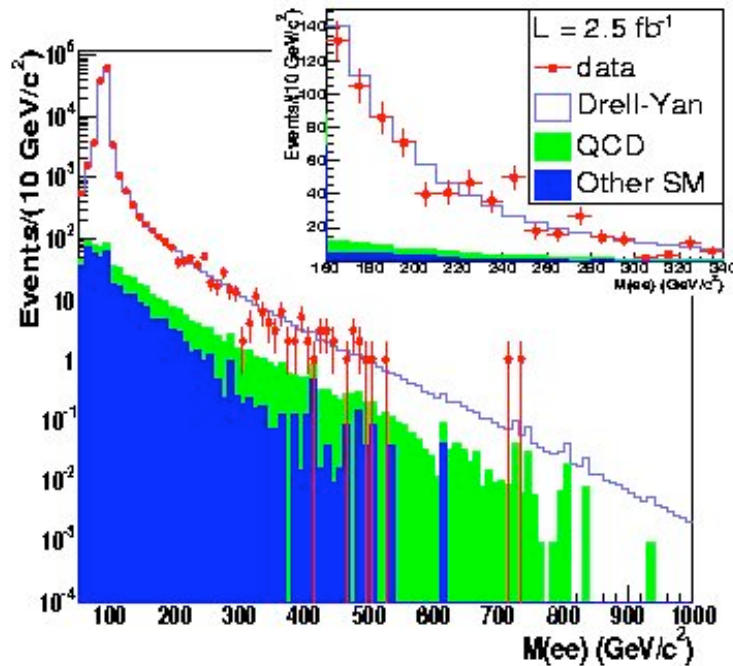


CDF Run II Preliminary

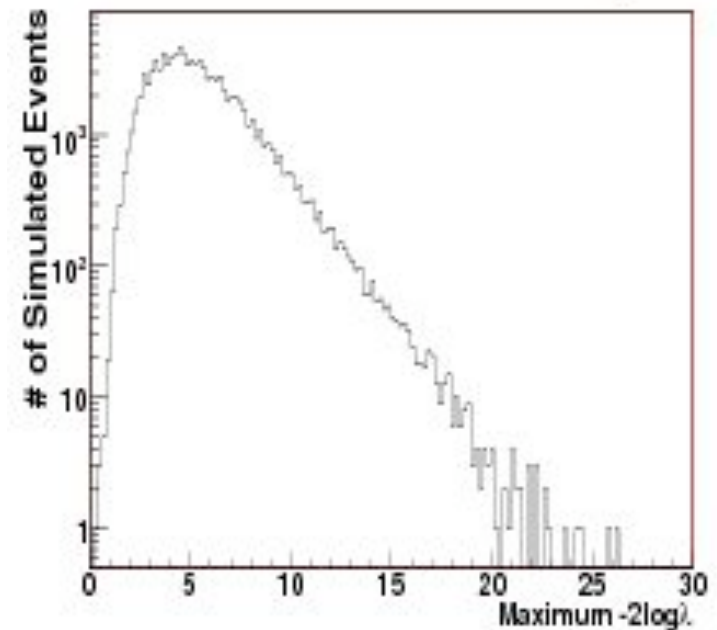


# Neutral Spin-1 Bosons: $Z'$

CDF Run II Preliminary



CDF Run II Preliminary

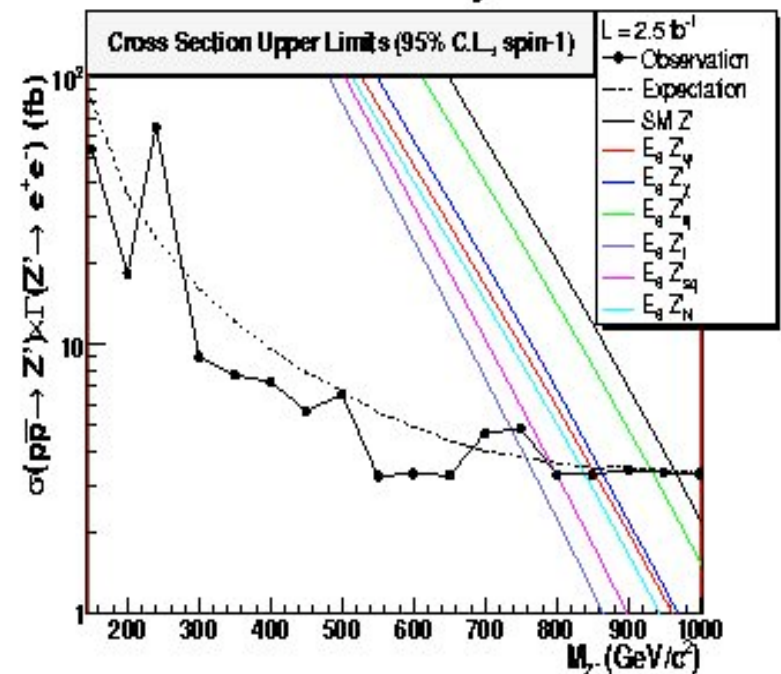


- $3.8\sigma$  excess at  $m(ee)=240 \text{ GeV}/c^2$
- But need to evaluate how likely this happens somewhere in the analysis!
- Technique:
  - Run Pseudo-experiments to see how often statistically one gets fluctuation of  $>3.8$
- Result: 0.6% probability ( $2.5 \sigma$ )
  - ♣ Need to see how new data look like

# Interpreting the Mass plots

- No evidence for any deviation from Standard Model => Set limits on new physics
  - Set limits on cross section x branching ratio
    - This is model independent, i.e. really what we measure
    - Any theorist can overlay their favorite curve
    - It remains valid independent of changes in theory
    - Always publish this!
    - Here:  $\sigma < 4$  fb at high mass
  - Can also set limits on  $Z'$  mass within certain models
    - This is model dependent
    - Nice though for comparing experiments, e.g. LEP vs Tevatron

CDF Run II Preliminary



For SM couplings:

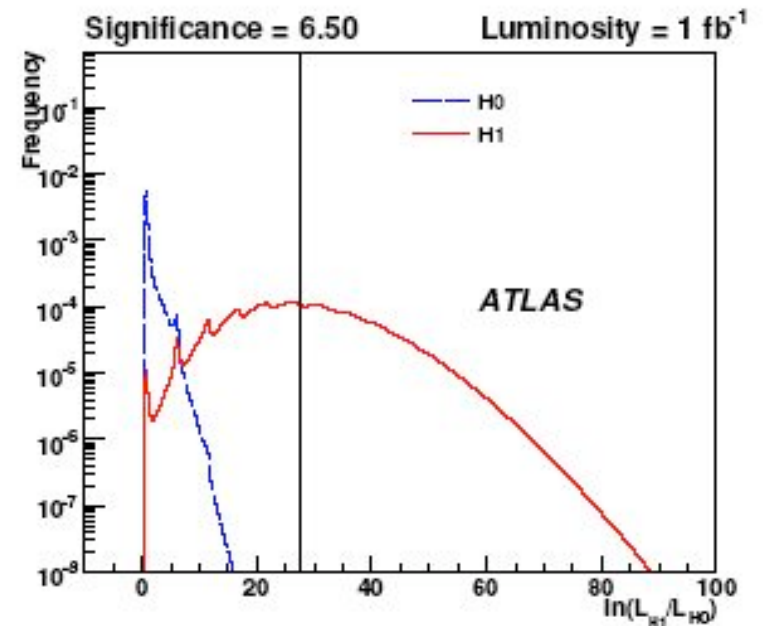
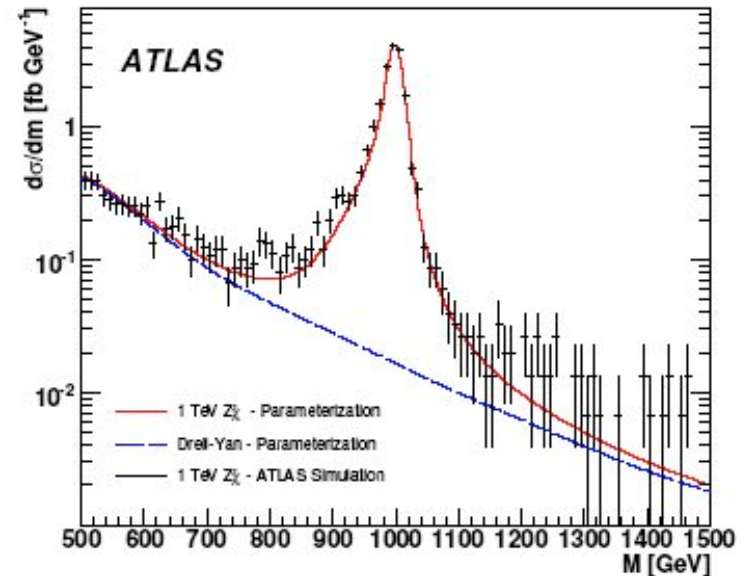
	$Z' \rightarrow ee$
limit	$M > 966 \text{ GeV}$

# Alternative Method

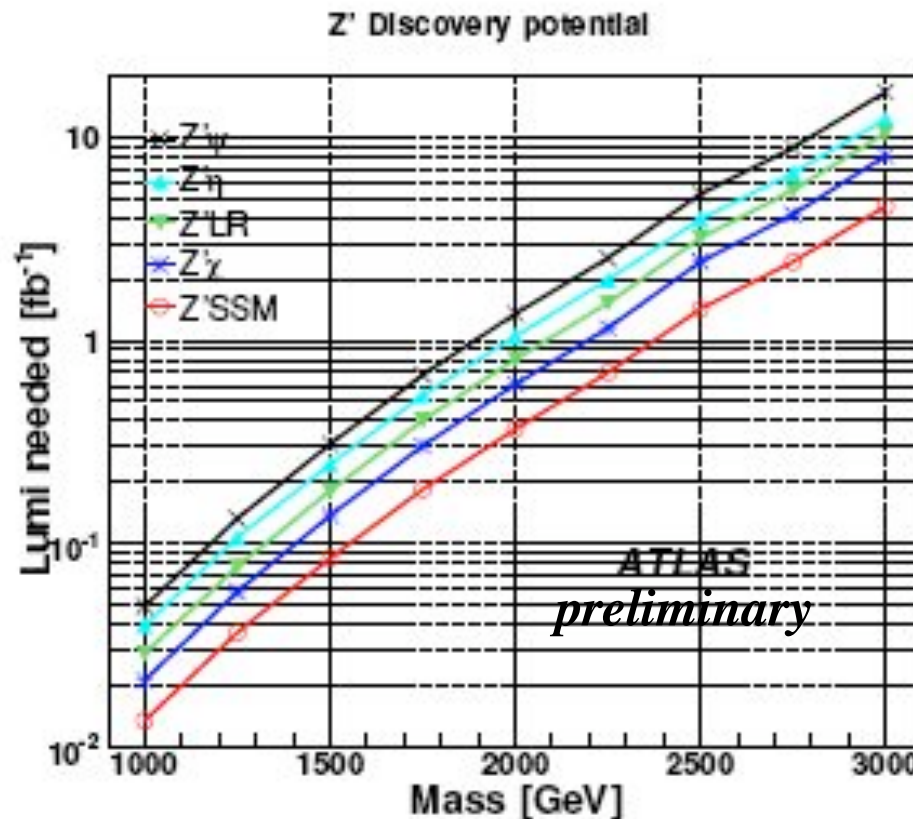
- Fit for bump in mass spectrum:
  - Parameterize distribution
  - Convolute with efficiency and detector resolution

$$\begin{aligned}
 \left. \frac{d\sigma}{dm} \right|_{\text{Signal}}(m) &= \frac{1}{m^2} \times G_{PDF}(m) \\
 &+ \mathcal{A}_{\text{peak}} \times \frac{\Gamma_{Z'}^2}{m_{Z'}^2} \frac{m^2}{(m^2 - m_{Z'}^2)^2 + m_{Z'}^2 \Gamma_{Z'}^2} \times G_{PDF}(m) \\
 &+ \mathcal{A}_{\text{interf}} \times \frac{\Gamma_{Z'}^2}{m_{Z'}^2} \frac{m^2 - m_{Z'}^2}{(m^2 - m_{Z'}^2)^2 + m_{Z'}^2 \Gamma_{Z'}^2} \times G_{PDF}(m)
 \end{aligned}$$

**Significance given by ratio of likelihood of signal-hypothesis vs likelihood of null-hypothesis**



# LHC $Z'$ Discovery Potential



- Probe  $>1$  TeV already with  $100 \text{ pb}^{-1}$
- Improves by  $\sim 1$  TeV for every factor 10 in luminosity
- Systematic uncertainties relatively negligible in these analyses